

モンテカルロ木探索を用いた統語情報を考慮した文生成

熊谷 香織[†] 持橋 大地[¶] 小林 一郎[†] 麻生 英樹[‡] Muhammad Attamimi[§]

中村友昭[§] 長井隆行[§]

[†]お茶の水女子大学大学院 [¶]統計数理研究所 [‡]産業総合技術研究所 [§]電気通信大学大学院

[†]{g1120515,koba}@is.ocha.ac.jp, [¶]daichi@ism.ac.jp, [‡]h.asoh@aist.ne.jp,
[§]m.att@apple.ee.uec.ac.jp, [§]tnakamura@uec.ac.jp, [§]tnagai@ee.uec.ac.jp

1 はじめに

自然言語処理の分野において、非言語情報を言語で表現する、文生成の研究が盛んになってきている。[1, 2, 3] 一方で、文生成の手法については、他の自然言語処理の手法に較べて研究報告が多くない。近年の一般的な文生成手法は、 n グラムモデル [4] やニューラル言語モデル [1] などの言語モデルの特徴のみを考慮したものが大勢を占めており、単語の依存関係など統語情報を考慮したものは数少ない。

このことを踏まえて、本研究では、統語構造を考慮した文生成手法を提案する。

2 モンテカルロ木探索を用いた文生成

2.1 概要

本研究では、文生成の手法として統語規則に文脈自由文法 (以下, CFG) を使用し、モンテカルロ木探索による適切な構造を持つ構文木の形成に基づき文を生成する。モンテカルロ木探索による統語構造の探索においては、構文木の生成において逐次的に CFG の文法規則を適用する。このとき、CFG による文法規則の局所的な関係しか見ることが出来ない。そのことから CFG に加えて、部分木に基づく構造的な素性を特徴量としたロジスティック回帰による正文 / 非文の分類を行うことで構造的な正しさを評価する。また、統合的な正しさの評価に加えて、同時に単語と単語の繋がりの良さの評価も考慮し、 n グラム言語モデルに基づく文のパープレキシティを用いた。上記二つの評価指標に基づき、適切な構文木の探索・形成を行うことで、より自然な文を生成できること示す。

2.2 モンテカルロ木探索

モンテカルロ木探索 (以下, MCTS) は、コンピュータ囲碁におけるゲーム AI の手法として注目されたが、ゲーム固有の知識を必要としないため様々なテーマに応用できる。MCTS の特徴として、ゲームの途中の状況を評価関数などを元に評価するのではなく、ゲームが終わるまで何度もシミュレーションを行い、その勝敗と探索回数を元に評価する。つまり、途中の段階を評価することが難しいタスクに適している。

構文木生成による文生成問題においても、生成途中の状態を評価することが難しいという課題がある。ことから、MCTS が文生成問題に適している。

MCTS のもうひとつの特徴として、Multi-Armed Bandit 問題に対処する指標として、Auer ら [6] によって提案された UCB (Upper Confidence Bounds) 1 値がある。UCB1 値スロットマシンを選択する指標として、従来の報酬の平均の代わりに用いられている。UCB1 値は、勝率の項、および探索が不十分なノードに対して選択の可能性を考慮した項から構成される (式 1)。

$$v_i + C \sqrt{\frac{\log N}{n}} \quad (1)$$

v_i はそのノードの勝率、 C は調整係数、 N は全試行回数、 n はそのノードを選択した回数を示す。UCB1 値における第 1 項が「知識の適用 (exploitation)」を、第 2 項が「探索 (exploration)」を考慮している。それによりバランスをとった探索が実行される。

2.3 MCTS を用いた文生成の処理の流れ

以下に、CFG を用いた文生成アルゴリズムとして、MCTS を適用した処理の流れを示す。

- step0. (初期設定): ルートノードに文の開始記号 S が適用される .
- step1. (選択): ルートノードから適用可能な文法規則を UCB1 値に基づいて選択する .
- step2. (拡張): 新たなノードを生成する .
- step3. (シミュレーション): 生成されたノードから文法規則をランダムに適用し終端記号の文字列を生成する .
- step4. (逆伝搬): 生成された文の文としての評価スコアが、他候補ノードの最大スコアを上回った場合 1 の値を (下回った場合は 0 の値を) 辿ってきた全てのノードに返し、勝率を更新する .
- step5. (ルートノードの更新): step1 から step4 を規定回数繰り返した後、ルートノードの子ノードの内、探索回数が最大のノードを次のルートノードとして更新し、step1 へ戻る .

3 文の評価

2.3 節の step4 において、文としての評価スコアという記述があるが、このスコアをどのように決めるかという問題が MCTS を用いた文生成において大きな課題となる . 囲碁においてゲーム終了時点での勝敗は明らかであるが、文生成においては生成された文が正文か非文かは明確な二値の状態としては与えられない . このことから本研究では、文の正誤の判断に 2 つの視点による評価値を導入した . 1 つは構文としての構造の正しさに対する視点、もう一つは単語の繋がりの良さに対する視点である . 以下、それぞれについて説明する .

3.1 構造の正しさに対する評価

構造の正しさを評価するため、構文木の部分木を素性としたロジスティック回帰による識別を行った . 岡野原ら [7] は、セマルコフクラスモデルによって単語の系列を意味的なクラスに分割し、そのクラスの N-gram を素性とし、文を正文と非文とに識別している . 本研究では、部分木を素性として採用し、正文・非文を判別する識別器を学習した . 作成した識別器により、生成された構文木が適切な構造から成り立っているかを構文木を構成する部分木の組み合わせより判断し、文が統語的に正しいか誤っているかを識別する .

学習データはコーパス中の部分木を正例とし、コーパスから得られた CFG を用いてランダムシミュレ-

ーションにより得られた構文木を疑似負例とした . この際、部分木抽出手法として freqt¹ を使用した . また、素性とする部分木は、構造のみに着目するために終端記号を除いたものとした .

3.2 単語の繋がりの良さに対する評価

単語の繋がりの良さを評価するため、直前の単語の種類数を重視するスムージング方法である、Kneyser-Ney スムージング [5] による 3 グラムのスコアを使用した . 以下に、3 グラムの場合の式を示す (式 2) .

$$p(z|xy) = \frac{(c(xyz) - d)}{c(xy*)} + \frac{d|xy*|}{c(xy*)} p(z|y) \quad (2)$$

c はカウント数、 d は種類数である . $d|xy*|$ が直前の単語の種類数を示す .

4 生成内容に対する制約の設定

前章までは、文生成における「いかに伝えるか (how to say)」について述べたが、「なにを伝えるか (what to say)」を表現することについても大きな課題である . 本研究では「なにを」表現したいかという情報が外部から得られたと仮定し、その情報をどのように文生成過程に組み込むかを考える .

例として、「dog」が「run」しているという状況を説明する .

「dog」の品詞は名詞である . いま、名詞の生成規則を、NN dog, NNS dogs のみに制限する . また、「run」の品詞は動詞であり、動詞の生成規則を、VB run, VBD ran, VBG running, VBP run, VBN run, VBZ runs のみに制限する . また、MCTS による文生成のシミュレーションにおいて、生成文中に名詞として「dog, dogs」が含まれていない場合、また動詞として「run, ran, running, runs」が含まれていない場合はその時点で負けとした .

以上のような制限を与えた範囲の中で MCTS による探索を行うことで、使いたい単語を使いつつ適切な構造と語順から成る文の生成を行なう .

また、使いたい単語の周辺の単語 (例えば時間情報や場所情報など) は、対象となるドメインにおいて適切に補われるべきである . 今回はそれらの情報は与えないため、使いたい単語を使用し、様々な状況を説明するような多様な文が生成できることを想定している .

¹<http://chasen.org/taku/software/freqt/>

5 実験

5.1 実験設定

コーパスは、Brown コーパス²を使用し、その中で文長が3~7の4,661文とした。文長を制限した理由として、生成内容として与える制約単語数から文生成に必要な文法は、短い文に出現する簡単な文法のみで十分と考えたからである。

コーパスをStanford parser³で構文解析し、解析結果をもとにCFGを作成した。CFGのサイズは、文法数7,220個で終端記号は5,867個であった。

次に構造的正しさを評価するための識別器の学習データを作成した。正例として、コーパスの4,661文をStanford parserで構文解析した構文木を使用した。

負例は、岡野原ら[7]の手法を参考に、コーパスより作成したCFGに基づいてランダムに生成した46,610個の構文木とした。負例数を正例の10倍とした理由としては、作成したCFGから生成され得る構文木の中で負例の方が圧倒的に多く存在することを考慮したためである。

また、MCTSのシミュレーション回数は10,000回に設定した。

5.2 統語構造に対する識別器の精度評価

作成したロジスティック回帰による識別器が正文と非文(疑似負例)を正確に識別できているかを確認するため10分割交差検定を行った。その結果98%の精度を確認した。

ここで作成した識別器のことを、以下「構造評価識別器」と称する。

5.3 構造的正しさの評価

まずは構造的な正しさのみに対し評価をする文生成実験を行った。2.1節のstep4における文としての評価スコアに、構造評価識別器により得られる確率を用いた。MCTSにより得られた生成文例を表1に示す。

上記の生成文を見るとSVOやSVの構造をした、構造的に正しい文が生成できていることが分かる。スコアをみるとおよそ0.99くらいであり、構造評価識別器により明らかに正例と判断される構文木が生成できていることも分かった。

²<http://clu.uni.no/icame/browneks.html>

³<http://nlp.stanford.edu/software/lex-parser.shtml>

表 1: 構造的正しさの評価をした場合の生成文例

生成文	構造識別器による確率
all mass nudged no teacher	0.999
this principle observed all super-condamine	0.999
all kay sank all round	0.999
some camping departs	0.994
those rim made these amount	0.999

一方で、単語の繋がりについては何も評価していないため、“some camping” “those rim made” など、不自然な単語の繋がりが見られた。また、語彙選択についても考えていないため、意味的にも不自然な文が生成された。

5.4 n グラムによる評価の追加

構造の評価に加えて、単語の繋がりについての評価もした文生成を行った。文のスコアを以下の様に設定した(式3)。

$$\begin{aligned} \text{文の評価スコア} = & \\ & \log(\text{構造評価識別器による確率}) \quad (3) \\ & + \log(1/3 \text{ グラムによるパープレキシティ}) \end{aligned}$$

また、上式の様に単語の繋がりとしてnグラムの評価を追加すると、短文がどうしても評価が高くなってしまふことから文長の制約を導入した。今回は、文長6未満の時はその時点で負けとした。

この時の生成文は表2の通りである。

表 2: n グラムによる評価を追加した時の生成文例

生成文	構造識別器による確率	3 グラムパープレキシティ
all of this may now ?	0.262	42.0
the innovation is a it ?	0.953	84.7
the the with would not ?	0.262	74.6
the lunch is what was a	0.996	120.5
both of the would not ?	0.262	56.8

以上の結果より、単語の繋がりの方については、“all of this” や、“both of the” などの慣用句的な表現が見られ、nグラムのスコアが影響していると考えられる。

構造的な正しさについてはおよそ0.9くらいになるときもあれば、0.25くらいになるときもある。すなわち、構造的に正しい文を生成できることもあれば、

できない時もある。n グラムによるスコアが追加されたことにより、構造的な正しさの評価の影響が小さくなってしまったと考えられる。

また、不自然な構文構造が選択された場合、単語の繋がりも不自然になることも確認された。

5.5 生成内容の設定

生成すべき内容を表す単語が得られたと仮定して文を生成することを考える。いま、「dog」が「run」している状況を説明したいとして、4 節で説明した生成内容に対する制約に基づき実験を行った。今回は単語数が 3~7 の文長を持つとして文生成を行った。まず、最も短い文の生成として文長 3 の制約を与えた。また、ある程度文を長くした際「dog」や「run」以外の単語を適切に補えるか検証するため、文長 5 と 7 で制約を与えた。文長制約の与え方としては、設定した長さよりも短かった場合はその時点で負けとした。また、勝敗の返し方は 5.3 節で示した方法を用いた。この時の生成文は表 3 の通りである。

表 3: 生成内容を仮定した場合の生成文例

文長	生成文	構造識別器 による確率	1 / 3 グラムパープレキシティ
3	the dog ran	0.994	793.7
5	the dog ran and in	0.925	423.7
5	the dog to run to	0.925	423.7
7	the dog running and to run to	0.767	421.9
7	why above comparative run to least dogs	0.00888	5291.0

以上のように文長が短い時は意味が通る文は生成できたといえる。文長が 7 など長くなってくると、名詞が dog または dogs だけだと語彙が少なすぎて意味が通らない文になってしまう。文長に合わせて、場所情報などを増やし、語彙を豊富にする必要がある。

6 おわりに

統語情報を考慮した文生成を行うため、CFG を適用規則とする MCTS により正文となる構文木の探索を行った。MCTS において、構文構造の正しさと単語の繋がり goodness に対する 2 つの視点における評価を行いながら探索を行うことにより、より自然な文の生成を目指した。

今回、2 つの視点に対する評価を行なったが、双方の評価値が安定して高くなるような探索の方法が確認できなかった。今後は、両評価値を共に効率良く収束するような探索の方法を考えるつもりである。また、表現したい内容を表した文の生成を目指すために、条件を様々に変更して検証するつもりである。

以上の点における課題を達成すれば、本研究は文生成を必要とする様々な自然言語処理の研究において有用なモデルになると考えられる。

参考文献

- [1] K.Xu, J.Ba, R.Kiros, K.Cho, A.Courville, R.Salakhutdinov, R.Zemel, Y.Bengio, Show, Attend and Tell: A Neural Image Caption Generation with Visual ,arXiv:1502.03044 [cs.LG], 2015.
- [2] M.Regneri, M.Rohrbach, D. Wetzell, S. Thater, B. Schiele, and M. Pinkal, Grounding Action Descriptions in Videos, Transactions of the Association for Computational Linguistics (TACL), 2013.
- [3] Haonan Yu and Jeffrey Mark Siskind, Grounded Language Learning from Video Described with Sentences, Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, pages 53-63, Sofia, Bulgaria, August 4-9 2013.
- [4] Yoshitaka Ushiku, Tatsuya Harada, and Yasuo Kuniyoshi. A Understanding Images with Natural Sentences. the 19th Annual ACM International Conference on Multimedia (ACMMM 2011), pp.679-682, 2011.
- [5] R. Kneser and H.Ney. Improved backing-off for m-gram language modeling. In Proceedings of ICASSP, Vol. 1, pp.181-184, 1995.
- [6] P.Auer, N.Cesa-Bianchi, and P.Fischer, Finite-time analysis of the multi-armed bandit problem, Machine Learning, 47:235-256, 2002.
- [7] D.Okanohara, and J.Tsujii, A discriminative language model with pseudo-negative samples, In Proceedings of ACL, 73-80, 2007.
- [8] Monte Carlo Tree Search (MCTS) research hub, <http://mcts.ai/>