2D5-OS-03b-3

動画像中の人の動作を表現する確率的言語生成に関する取組み

An Approach to Probabilistic Text Generation of Human Motions in Video Clips

*¹お茶の水女子大学大学院 人間文化創成科学研究科 理学専攻 情報科学コース Advanced Sciences, Graduate School of Humanities and Sciences, Ochanomizu University

*2独立行政法人産業技術総合研究所知能システム研究部門

Intelligent Systems Research Institute, National Institute of Advanced Industrial Science and Technology

In this study, we propose a framework for probabilistic text generation of visual information. As for visual information, we particularly employ the time-series data of human motions captured by a Kinect camera. The time-series data are applied by several dimension reduction procedures and then turned to be the form which can be applied to machine learning. A pair of the analyzed time-series data and its intermediate representation which corresponds to the semantics of the human motion is learned by a log-linear model. As linguistic resource to generate a text, we collected various natural language expressions for human motions and build a bi-gram model for each motion. In our framework, once the intermediate representation is decided by observing time-series data; a proper bi-gram model corresponding to the intermediate representation is chosen; and then a text is generated by solving dynamic programming of the bi-gram model. Through experiments to generate texts describing human motions, we have confirmed that our proposed framework works well.

1. はじめに

近年、大量の動画像データを取得することが容易になってきている。一方で、大量に収集したデータを有効に利活用出来ているとは言えない。例えば、監視カメラの動画像データに映る内容を把握するためには、全てを人目で見る必要があるが、データの多さに応じた時間を要してしまう。もし、大量の動画像データから特徴的なイベントを捉え、またそのイベントを言葉として表現することが出来たら、動画像データに映る内容を簡単に把握できるとともに、言葉で動画像中のイベントの検索も行うことができると考える。そこで本研究では、動画像の情報を入力とした確率的なテキスト生成手法を提案する。

2. 研究概要

本研究の概要を図1に示す.まず, Kinect*1がもつ人の骨格を追跡するライブラリを用いることで, 人の動きを時系列データとして取得する.取得された時系列データはいくつかの次元圧縮作業を行い, データと自然言語をつなぐ中間表現とともにデータベースに格納される. その後, データベース内に蓄積された時系列データと中間表現の対応関係を機械学習することで,動作判別器を生成する.テキスト生成に用いられる言語資源は,人の動作の表現を被験者実験によって収集し,それぞれの中間表現に対してバイグラムモデルを構築する.これにより中間表現を選択すると,その中間表現に対応したバイグラムモデルが選択され,そのモデルに動的計画法を適用することで,人の動作を表現するもっともらしい語の組み合わせを選ぶことができる.

連絡先: 小林瑞季,お茶の水女子大学大学院人間文化創成科学研究科理学専攻情報科学コース 小林研究室,〒 112-8610 東京都文京区大塚 2-1-1,03-5978-5708,kobayashi.mizuki@is.ocha.ac.jp

*1 http://www.microsoft.com/en-us/kinectforwindows/

3. 関連研究

マルチメディア情報を入力としてテキスト生成を行う関連研 究として、Barbuら [Barbu 2012] のショートビデオの説明文を 生成するシステムが挙げられる. これらの説明文は、誰が何に何 をしたのか、どこでどうやって行ったのかを説明している. これ は、簡単な文法を加味したテンプレートベースのテキスト生成 を行っている. また、柔軟なテキスト生成に関しては、Belz と Kow ら [Belz 2007, Belz 2009] の生成空間の統括的なモデルを 用いた確率的生成手法が挙げられる. 本研究により関連の深い研 究として, Liang ら [Liang 2009] や Angeli ら [Angeli 2010], Konstas ら [Konstas 2012a, Konstas 2012b] の研究が挙げら れる. Liang ら [Liang 2009] は、テキストと意味との関係を学 習する手法を提案しており、そこでは、イベントはデータベー スのレコードで表せると仮定し、レコードと自然言語で表記さ れた説明文との関連を機械学習によって取得する. Angeli ら [Angeli 2010] は、Liang ら [Liang 2009] が提案したモデルに 基づく潜在情報と表層情報をテキスト生成する手法を提案して いる。また、Konstas ら [Konstas 2012a, Konstas 2012b] は、 入力情報固有の構造を説明する確率的な自由文脈文法を定義し ており、Liangら [Liang 2009] や Angeliら [Angeli 2010] と 同様に、データベースのレコードと説明文を用いている. 彼ら は、重みを加えたグラフによって文法を表現し、また与えられ た入力に対しもっとも適切な導出木を見つけることでテキスト 生成を行う

本研究では、これらの関連の高い研究の考え方を参考にし、 時系列データと動作を表す自然言語の説明文との対応を学習させるために対数線形モデルを採用した.提案するテキスト生成の手法は単純であり、生成に文法を必要とするような複雑な文は生成することができないが、動画像を入力とし、素朴でもっともらしい文を容易に生成することが可能である.

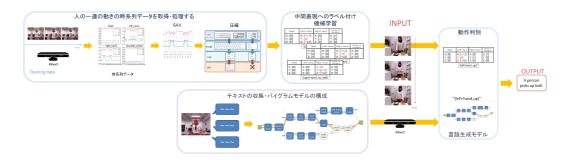


図 1: 動画像を入力とする確率的テキスト生成の枠組み

3.1 時系列データ処理

人間の動作の時系列データは、Kinect カメラを用いて取得する。Kinect の開発元である MicroSoft 社は、人間の骨格を推定できる標準ライブラリも提供しており、そのライブラリを用いると人の関節の 3 次元情報を推定することができる。本研究では、RGB 画像と深度センサー、またそれらを用いた人物の関節位置推定も用い、RGB 動画像と人物の頭・肩の中心・右手・左手の 4 箇所の xyz 座標の時系列データを取得する (図 2 参照)。

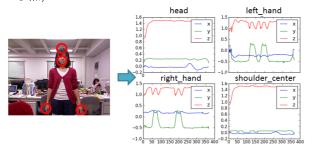


図 2: Kinect を用いた時系列データ取得

人の骨格を追跡することで得られた時系列データは、Symbolic Aggregation approXimation (SAX) [Lin 2003] を使い、文字列に変換する.

SAX によって変換して得られた文字列から動作とみられる 個所を取り出す. ここでは、ある動画像データ中の全ての文字 列において一つ前の文字から変化がなければ「動きがない」、変化があれば「動きがある」とみなす (図 3 参照).

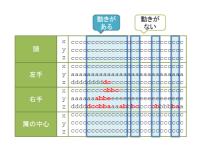


図 3: 動きの抽出例

その後、「動きがある」とみなされた個所の文字列を変化量(図4中のアルファベットの下の数値)に変換し、圧縮する(図5参照). これは同じ動作でも位置やスピードによっては文字列がある一定の間隔でずれたり文字列の長さが変化したりしてしまい、同じ動きとして学習されないためである. これにより、一定の間隔でずれてしまったものも長さが違うものでも、同じ動きとしてとらえることを可能とする. また、より特徴的

な動作を抽出するために、圧縮された変化量うち最大の大きさが2未満を示す動きは取り除く(図5参照).

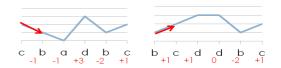


図 4: 文字列の変化量

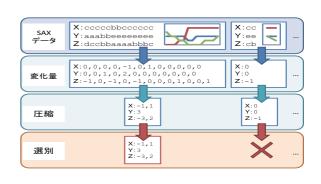


図 5: データの圧縮・選別の例

3.2 中間表現

テキスト生成では、時系列データと自然言語文をつなぐ中間 表現を用いることでテキスト生成に使う言語資源を選択する. 中間表現は表1のように定義する.

表 1: 中間表現 action 中間表現 'up{object}" upward movement down "down{object} ment left "to_left{object} leftward movement "to_right{object}" "pass{object1,object2} rightward movement right cooperative move-

3.3 時系列データの動作判別

本研究では人の動作の判別を行うために対数線形モデルを用い、処理された時系列データと中間表現の対応を機械学習させる。3.1 で述べた時系列データ処理を施したデータ d と、人の動作を表す中間表現 r から構成した素性ベクトル ϕ を用いて、式 (1) の対数線形モデルを構成することで、データが与え

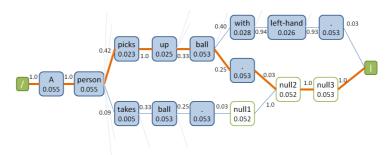


図 6: null ラベルを用いたバイグラムのイメージ

られた下での各言語表現が選ばれる確率 P(r|d) をモデル化し た. ここで、 $Z_{d,\mathbf{w}}$ は正規化係数である.

$$P(r|d) = \frac{1}{Z_{d,\mathbf{w}}} \exp(\mathbf{w} \cdot \phi(d,r))$$
 (1)

3.4 バイグラムモデルによるテキスト生成

本研究では、バイグラムモデルを用いた単純なテキスト生成 を行う. それぞれの動作に対しバイグラムモデルを構築するた めに被験者実験を行い、特定の動作に対して様々な自然言語表 現を集めた.これにより、観測された時系列データに対して特 定の中間表現が与えられたとき, 言語資源としてバイグラムモ デルを選択しテキスト生成を行う. しかし, 例えば同じ動作で も, ある人は10語で表現し, またある人は15語で表現するな ど、複数の表現方法がある. このことから、文の長さに依存し ないテキスト生成が行えるよう,バイグラムモデルに null ラ ベルを導入した. null ラベルは、文の中の単語として扱われ、 他の単語と同じようにユニグラムとバイグラムをとられる. こ のように null ラベルを扱うために、動的計画法を適用する前 に以下に続く前処理をそれぞれの文に対して行う. まず, 全て の文で単語数の最大値 max, 最小値 min を得る. 次に, max から min を引き, null に振る番号の最大値 null_max を求め る. 最後に、それぞれの文に対し、単語数が max に満たなけ れば、null_max から1ずつ引いた値を、足りない数だけ文末 から文頭に向け挿入していく. null ラベル導入のイメージを, 図7に示す.

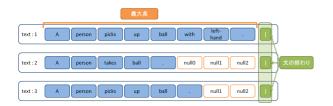


図 7: null ラベル導入のイメージ

文中の各 null ラベルに違う番号をつけることによって別の 単語として見なし、それぞれがバイグラムモデルの1要素とし て扱う. また、本研究ではバイグラムモデルを構築する際に、 使用する文の取捨選択を行わないことで,多くの語と関連づけ ることができるため、より複雑なテキスト生成を行うことがで きる. 人の動作「pick up ball」を説明したバイグラムモデル のイメージを図6に示す.

人の動作を説明するもっともらしい文を生成するためには, このバイグラムモデルに動的計画法を適用することで得られる.

実験

ここでは,「ボールを持ち上げて箱に入れる」という簡単な 動作(図8)を言葉で表現することを目的とする.



図 8: 言語化の対象となる動作

実験仕様

まず、言語化の対象となる動作を「pick」「pass」「put」の 3つの基本動作から成ると定義しておく. これは, 動作のどの 部分を自然言語文で説明するのかを自動で決定することが難し いためである. ここでは、それぞれの動作に対し、自然言語で の説明文を生成することとする. 被験者実験として, 対象とな る人の動作の Kinect ビデオを観賞し、それについて自然言語 で説明してもらうという実験を12人に対し行った.収集した 日本語の説明文を英訳し、これを言語資源としてバイグラムモ デルを構築した. 言語資源となった英文の全文数, 語数, 語の 種類数を表2に示す.

動き 語の種類数 文数 語数 1 33 214 47 2 148 28 18 3 36 290 28

表 2: 収集された文の特徴

動作判別には3.3で示した対数線形モデルを適用し、テキス ト生成に使われる中間表現の判別に用いた. この識別器は, 対 象となる動作を捉えた 20 のデータを 15 の訓練データと 5 の 評価データに分割し5クロスバリエーションした結果、精度 の平均が84%を示した.

4.2 実験結果

中間表現は、1つ目の動作は "up ${object}$ " 、2つ目の動作は "pass{object1,object2}", 3つ目の動作は"down{object}"" と認識された. 次に、選ばれた中間表現に対してあらかじめ構 築されたバイグラムモデルに動的計画法を適用することで、動 作を説明するもっともらしい文を生成する.

結果として、それぞれの動作に対して尤度の高かった上位3 文を表3に示す.

表 3: 各動作に対する生成文の上位 3 文

20 0 1 20 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2		
動作	生成文	尤度
	• A, person, picks, up, pink, ball, ., null_8, null_9, null_10, null_11, null_12, null_13, null_14, null_15,	5.68e-24
1	• A, person, picks, up, ball, with, left-hand, ., null_8, null_9, null_10, null_11, null_12, null_13, null_14, null_15	2.52e-24
	• A, person, picks, up, pink, ball, with, left-hand, ., null-8, null-9, null-10, null-11, null-12, null-13, null-14	2.10e-24
2	• A, person, passes, ball, to, right-hand, ., null_7, null_8, null_9, null_10, null_11	6.29e-16
	• A, person, passes, red, ball, to, right-hand, ., null_7, null_8, null_9,null_11, null_10	3.08e-18
	• A, person, passes, ball, from, left, to, right-hand, . , null_7, null_8, null_9	2.05e-18
3	• A, person, puts, ball, in, box, ., null_8, null_9, null_10,	4.90e-15
	• A, person, puts, ball, in, box, ., null-7, null-8, null-9, null-10	1.22e-15
	• A, person, puts, ball, to, another, box, ., null_7, null_8, null_9	2.16e-16

4.3 考察

実験結果から、人の動作を正確に表現する文が生成出来ていることが確認できた。また、表3の生成文をみると、いくつかの文で終端文字「|」が出てきていないことが分かる。これは、バイグラムモデルが集めた文に現れる語のバイグラムの組み合わせによって構成されているためである。これにより、バイグラムモデルへnullラベルを加えた文が、集められたどの文よりも長く生成される可能性がある。また一方で、文が長くなればなるほど、その文の尤度が低くなっていく。したがって、集められた文より長い文は生成されないという仮定の下で、集めた文の最大の単語数を生成文の単語数とした。

5. まとめと今後の課題

本研究では、動画像中の人の動作を表現する確率的言語生成の枠組みを提案した。Kinect ビデオで抽出された人の動作は、時系列データとして取得され、いくつかの次元圧縮手法を適用することで機械学習に適した形に変換される。また観測された人の動きを表現するために、被験者実験によって集められた自然言語文に基づきバイグラムモデルを構築し、動的計画法を適用することで、もっともらしい語の組み合わせを取得する。さらに、バイグラムモデルに番号を付けた null ラベルを導入することにより、文生成に単語数の制限をつけずに自然言語文生成を行うことができた。また、提案手法はテンプレートによるテキスト生成ではなく、確率的なモデルによる生成であることから、例えばさらに文を収集すればそれに合わせて出力文も変化していくなど、資源となる文書によって様々な自然言語表現を得ることができる。

一方で、現段階では構文制約や物体認識を取り入れてはいない。そのため今後の課題として、こうした知識を導入するとともに、より正確にイベントを説明するようなテキスト生成が行えるよう発展させていきたいと考える。また、中間表現とバイグラムモデルとの対応付けをより柔軟したり、一連の動作から自然言語文によって説明される動作を区切る問題にも取り組んでいきたい。

参考文献

- [Bangalore 2000] Srinivas Bangalore and Owen Rambow 2000. Exploiting a Probabilistic Hierarchical Model for Generation, Proceedings of the 18th conference on Computational Linguistics (Coling 2000), Volume 1, pp.42-48,
- [Barbu 2012] A. Barbu, A. Bridge, Z. Burchill, D. Coroian, S. Dickinson, S. Fidler, A. Michaux, S. Mussman, S. Narayanaswamy, D. Salvi, L. Schmidt, J. Shangguan, J. Siskind, J. Waggoner, S. Wang, J. Wei, Y. Yin, and

- Z. Zhang. 2012. Video In Sentences Out, Conference on Uncertainty in Artificial Intelligence (UAI),
- [Belz 2007] Anja Belz, 2007. Probabilistic Generation of Weather Forecast Texts, Proceedings of NAACL HLT 2007, pp.164-171
- [Belz 2009] Anja Belz and Eric Kow, 2009. System building cost vs. output quality in data-to-text generation, In Proceedings of the 12th European Workshop on Natural Language Generation (ENLG-09, pp.16-24, Athens, Greece
- [Angeli 2010] Angeli, Gabor and Liang, Percy and Klein, Dan, 2010. A simple domain-independent probabilistic approach to generation, Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pp. 502–512, Cambridge, Massachusetts
- [Konstas 2012a] Konstas, Ioannis and Lapata, Mirella, 2012. Unsupervised concept-to-text generation with hypergraphs, Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Montreal, Canada, pp.752-761
- [Konstas 2012b] Konstas, Ioannis and Lapata, Mirella, 2012. Concept-to-text generation via discriminative reranking, Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1, pp. 369–378, Jeju Island, Korea
- [Lin 2003] Lin, J., Keogh, E., Lonardi, S. and Chiu, B. 2003. A Symbolic Representation of Time Series, with Implications for Streaming Algorithms DMKD' 03
- [Lapata 2003] Mirella Lapata, 2003. Probabilistic Text Structuring: Experiments with Sentence Ordering, In Proc. of the Annual meeting of the Association for Computational Linguistics pp.545–552
- [Liang 2009] Percy Liang, Michael I. Jordan, Dan Klein 2009. Learning Semantic Correspondences with Less Supervision, ACL-IJCNLP
- [Lu 2009] Wei Lu and Hwee Tou Ng and Wee Sun Lee, 2009. Natural language generation with tree conditional random fields, EMNLP, pp.400–409