

記号化された運動と自然言語の相互連想モデルと コミュニケーションへの応用

Association between Symbolized Motion and Natural Language, and its Application to Communication

川辺直人 高野渉 中村仁彦
Naoto Kawabe Wataru Takano Yoshihiko Nakamura

東京大学大学院 情報理工学系研究科 知能機械情報学専攻

Department of Mechano-Informatics, Graduate School of Information Science and Technology, The University of Tokyo

Natural language is the symbol that represents the world. We have developed the symbolization of human and robot motions by Hidden Markov Models (HMMs), while researches in computational linguistics have demonstrated the statistic methods of natural language processing. In this paper, we present an association model that connects symbolized motions and natural language with Conditional Random Fields (CRFs). CRFs enable this model to interpret a motion in a proper context by referring the sequence of motion symbols. In addition, we propose its application to the decision of an action based on the Theory of Mind (ToM). A robot has self and other's models, and the recursive computation of an interaction helps it select the best motion to take. We integrate the association model and a model representing relationship between sentences to realize that computation. This research contributes to the construction of the human-robot interaction.

1. はじめに

人間とコミュニケーションをする際、ロボットによる言語獲得は有効かつ重要である。言語は、自然や運動を分類・記号化して扱う物である。[Inamura 04] は隠れマルコフモデル (Hidden Markov Model: HMM) を用いた運動の記号化を行った。さらに、自然言語処理に用いられる言語モデルを組み合わせることで、[Takano 08] は言語と運動の連想モデルを提案した。本研究では、言葉の選択に運動の文脈を考慮することのできる条件付確率場 (Conditional Random Field: CRF) を用いた連想モデルを提案する。また、このような言語の獲得は推論を可能とする一つのステップであると考えられる。特に他者の推論は人との共存の上で不可欠である。人間が行う他者の推論を工学的に実装した研究では、再帰性に注目した心の理論がある [伊藤 99]。他者の行動を言語上で再帰計算により推論することで、他者とのインタラクションを予想し、その結果から自分の行動を選択する手法を説明する。これは運動の言語化に基づいた他者の理解の一例である。

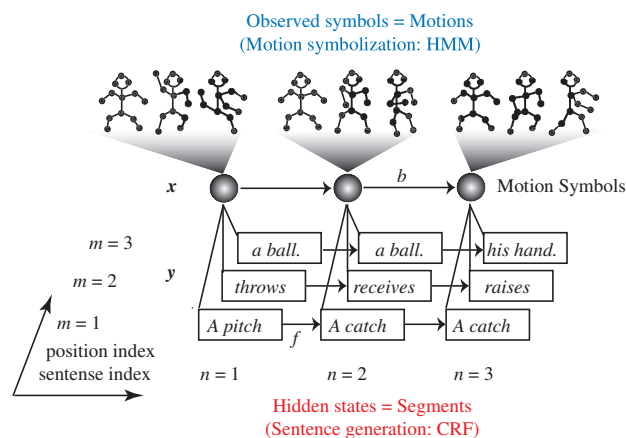


図 1: Motion-language associative model

2. 運動・言語相互連想モデル

2.1 運動・言語モデル

まず、図 1 に示した運動の言語化について述べる。このモデルでは、入力運動は HMM によって記号化されているとする。また、自然言語の文はあらかじめ構文解析機 CaboCha [Kudo 02] によって文節に分けられている。 n 番目の運動入力を x_n 、文中 m 番目の位置にある文節を y_{mn} とすると、運動記号列 x から対応する文節列 y_m を推定する確率は次式で表される。

$$P_{\eta}(y_m|x) = \frac{\exp\left(\sum_i \eta_i f_i(x, y_m)\right)}{\sum_{y \in Y_m} \exp\left(\sum_i \eta_i f_i(x, y)\right)} \quad (1)$$

上式中、 f は素性関数といい、運動記号と文節の変化に応じてバイナリを返す。 η_i はパラメータ、 Y_m は y_m の候補の集合である。素性関数 f は次のように設定した。

$$f_{mi} = \begin{cases} b_j & \text{if } y_{mn} = w_{mk}, y_{m(n+1)} = w_{ml} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

$$b_j = \begin{cases} 1 & \text{if } x_n = \lambda_p, x_{n+1} = \lambda_q \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

f_{mi} は文節 w の変化を、 b_j は運動記号 λ の変化を観測している。学習では式 (1) の対数尤度を取り、最急降下法によってパラメータ η_i の最適化を行う。また、学習と同時に運動記号と文節の対応関係を記した辞書を構築する。

実行時、運動記号列 x が与えられると、文中位置 m における最適な文節列 y_m は、次のステップに従い図 2 のように求められる。これを各位置 m で行い、得られた文節を並べることによって 1 文を出力する。

Step1 運動記号に対応する文節候補を辞書から選択する。

連絡先: 中村仁彦, 東京大学大学院情報理工学系研究科, 〒113-8656 東京都文京区本郷 7-3-1, Tel: 03-5841-6381, Fax: 03-3818-0835, E-mail: nakamura@ynl.t.u-tokyo.ac.jp

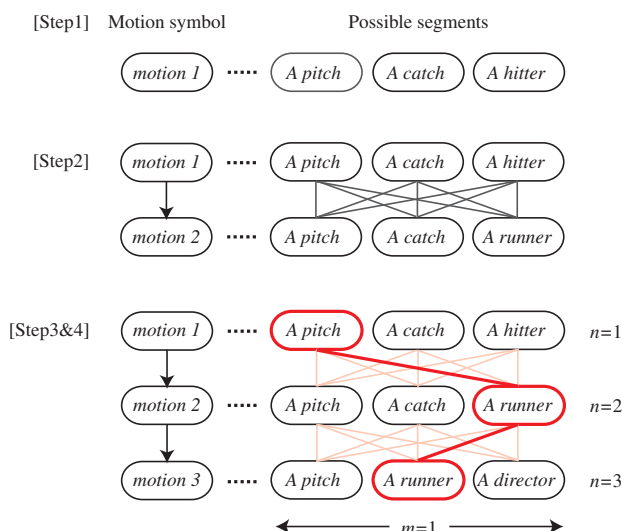


図 2: Viterbi search

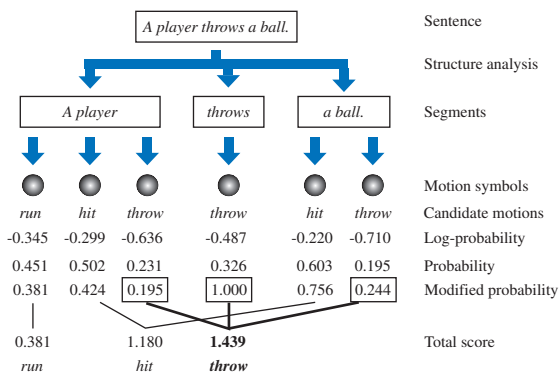


図 3: Motion symbol selection

Step2 n 番目, $n+1$ 番目の文節候補 $y_{mn}, y_{m(n+1)}$ に対し次式を計算する.

$$\sum_i \eta_i f_i(x_n, x_{n+1}, y_{mn}, y_{m(n+1)}) \quad (4)$$

Step3 $y_{m(n+1)}$ に接続する候補 y_{mn} の中で, 最大の評価値を持つパスとその時の値を記憶する.

Step4 n が終端に達したところで, 最大の評価値を持つ物を探し, 遡ることでパスを決定する.

2.2 逆連想

言語から運動を推定することをここでは逆連想と呼ぶ. 前項の確率式 (1) を用いると, ベイズの定理から

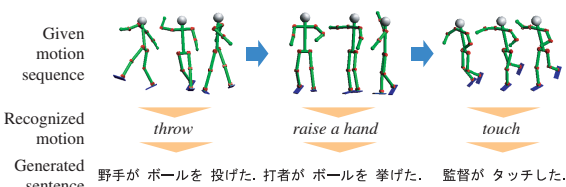
$$P(x|y_m) \propto P(x)P_\eta(y_m|x) \quad (5)$$

である. ここで, $P(x)$ は Bigram, すなわち,

$$P(x) = \frac{N(x_{n+1}, x_n)}{N(x_n)} \quad (6)$$

とする. $N(\cdot)$ は学習データ中に出現した要素をカウントする関数である. ある文が入力された時, それを CaboCha を用い

Result (A): Sentence generation



Result (B): Motion selection

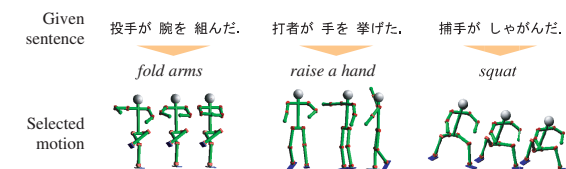


図 4: Experimental results

て構文解析し, 各文節から辞書を参考に候補となる運動記号を挙げ, 各々について式 (3) 右辺の対数をとった次式を用いて評価を行うことで運動記号の選択が可能となる.

$$\sum_n \log P(x_{n+1}|x_n) + \sum_n \sum_i \eta_i f_i(x_n, x_{n+1}, y_n, y_{n+1}) - \log \sum_{y \in Y^m} \exp \left(\sum_n \sum_i \eta_i f_i(x_n, x_{n+1}, y_n, y_{n+1}) \right) \quad (7)$$

出現頻度の低いユニークな運動を選択に反映させるために次の補正を行う. 図 3 のように, 各文節で得られた評価値を指数をとって確率に直し, 文節の候補の中でそれらの和が 1 になるように調節する. 最後に, 候補の運動について和をとり, 最も高い値を持つ物を選ぶ.

2.3 相互連想の実験

上記のモデルの有効性を確認するために, 野球に関する 17 個のモーションキャプチャデータを用いて実験を行った. 複数の運動を並べ, 対応する日本語文は手で与えた. 2 個を 17 セット, 3 個を 34 セット用意した. 式 (1) のパラメータ η を学習させ, これを用いて運動から文の連想, 逆連想を行った. 実験結果を図 4 に示す.

図 4(A) は運動の言語化を行ったもので「投げる」(throw) から「野手がボールを投げた」のように観測された運動記号から正しい文が選択されていることが分かる. また, CRF を用いたことにより, 同じ運動が前後の文脈によって違った解釈をされることも確認されている. 例えば「腕を挙げる」(raise a hand) 動作に対し, 図 4 では「打者がボールを挙げた」と解釈されているが, 前後が異なる別の運動記号の文脈では「腕を挙げる」のまま解釈されることが確認されている.

図 4(B) は, 言語からの運動生成を行った結果で「投手が腕を組んだ」という文から「腕を組む」(fold arms) という運動記号を選択している. 与えられた言語から運動の言語化で得られたパラメータを再利用して適切な運動を選択していることが分かる.

以上から, このモデルが有効にはたらいっていることが示された.

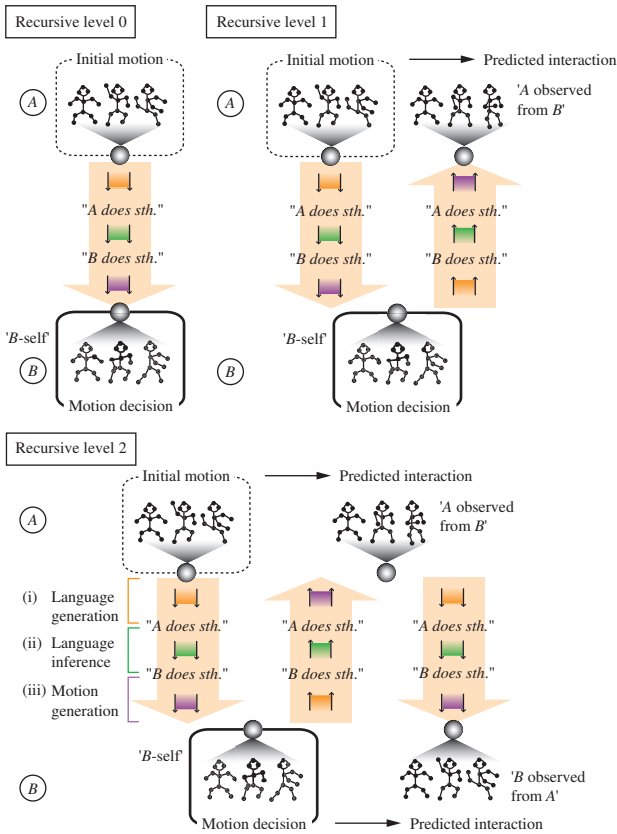


図 5: Theory of Mind

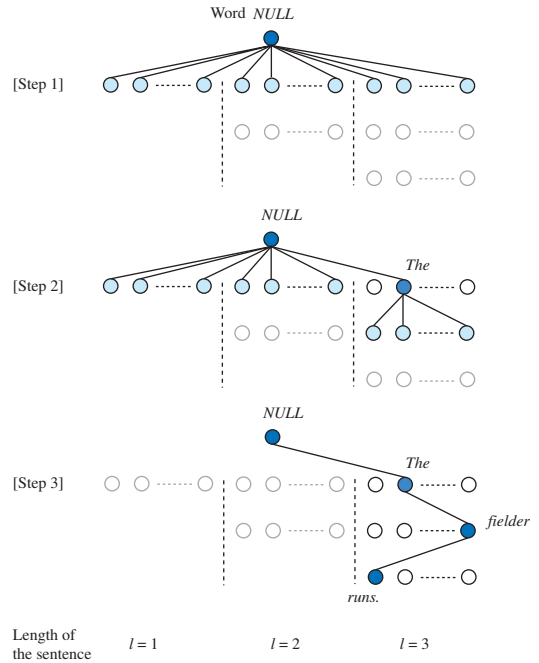


図 6: A* search

ここで, $t(f|e)$ は生成される前の文節 e から生成後の文節 f への生起確率, a は文長 m, l の被生成文, 生成文において j 番目の被生成文節が i 番目の訳語に接続する確率である. 2つのパラメータ t, a は EM アルゴリズムによって最適化される. ベイズの定理から,

$$P(e|f) \propto P(e)P(f|e) \quad (9)$$

となる. $P(e_i|e_{i-1})$ は言語モデルと呼ばれ, ここでは Bi-gram を設定した. すなわち,

$$P(e) = \prod_i P(e_i|e_{i-1}) \quad (10)$$

である. 以上を用いて与えられた文に対して式 (11) を評価値として図 6 に示すように A* サーチによって対応する文の検索を行う.

1. l が不定のため複数の文長を用意し, 各文長の先頭の文節 e_1 において次式を計算する.

$$\sum_i \log P(e_i|e_{i-1}) + \sum_j \log \sum_i t(f_j|e_i)a(i|j, m, l) \quad (11)$$

2. 上式で最大の評価値を持つ単語を選択し, その文長を持つ文において, 次の単語で式 (11) を計算する.
3. 2. の単語の最大値と, 1. の他の文長の単語の最大値を比較し, 大きい方の単語を選択する.
4. 以上を繰り返し, 選択された単語が文長に到達した時点で探索を終了する.

A が初期運動をとった時, B がインタラクションを考えるに当たって, 次の 2 点を考慮する.

1 点目は, 文の探索に用いるパラメータを立場によって複数用意する点である. 例えば, B は自分の反応に対する A の反

3. コミュニケーションへの応用

3.1 心の理論

次に, 運動の言語化の技術をベースとした他者のモデル化を行う. 言語を操ることで, 他者とのやりとりを経験から想像することが可能になる. この他者を自己の中に持つことを心の理論 [Premack 78] は明らかにしている. そこで, 心の理論の工学的応用を言語を介して実現し, 他者の理解について考える.

心の理論の定式化において, 再帰レベル (再帰の深さ = 何手裏 [先] を読むか) を設定する [Takano 05]. いま, コミュニケーションを行う 2 者 A, B を仮定し, そのやりとりを再帰計算によって行う. A の運動が観測された時, B はそれを言語化し, 表わされた文に最も相応しい対応を言語上で検索する. 検索には IBM 翻訳モデル 2 を適用する. この翻訳モデルは単語とその配置のみをパラメータとしたモデルで, 因果に当たる文のペアを教師として与えれば因果関係を獲得すると考えられる. 図 5 のように, 運動の言語化・対応する文の検索・文からの運動生成を 1 セットとする. 再帰レベルが 0 の時は相手の運動に対する自分の反応を考え, 再帰レベルが 1 になるとさらに自分のとった反応に対する相手の反応まで考慮する. 設定した再帰レベルまで繰り返し行うことで B は 2 者間のインタラクションを予想し, いま取るべき運動を決定する.

3.2 他者の運動の予測と報酬

他者の運動の予測には翻訳モデル [Brown 93] を用いる. 翻訳モデルでは式 (8) に示した翻訳確率を用いる.

$$P(f|e) = \prod_{j=1}^m \sum_{i=0}^l t(f_j|e_i)a(i|j, m, l) \quad (8)$$

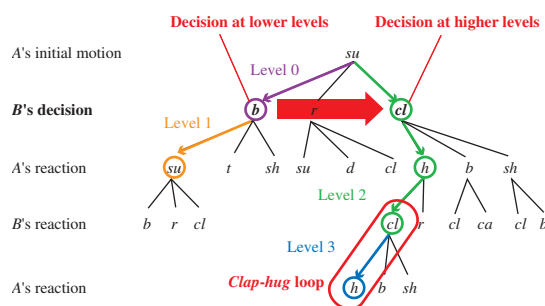


図 7: Experimental result of the theory of mind: *su* = stand up, *b* = bow, *r* = raise a hand, *cl* = clap, *t* = touch, *sh* = shake hands, *d* = dash, *h* = hug. The motion bow is selected at level 0 and 1. As the recursive level increases, the selection converges to *clap* because there is a *clap-hug* loop which attracts the interaction.

応を *B* から見た *A* として獲得する。同様に、*A* から見た *A*, *B* から見た *B*, そして *A* から見た *B* を獲得する。

2 点目は、現在の対応が目的のインタラクションにどれだけ近いかを表す指標として文の長さを報酬 $L(e)$ として用いる点である。コミュニケーションにおいては、ある対応した運動が起こりやすいというだけでなく、目的の状態に誘導することが考えられる。パートナーの様子を窺うには声色や表情といったものが考えられるが、ここでは簡単のため文長を報酬として与え、記述が長い運動へと選択を変えていく。例えば「投手がボールを投げた。」であれば、そのバイト数 22 が報酬となる。これを式 (11) にかけることで運動選択の際の評価の値とする。

$$P(f|e) \cdot L(e) \quad (12)$$

3.3 シミュレーション

再帰計算による心の理論の計算をシミュレーションにより行う。*A* の運動が与えられた時に、*B* がどのようなインタラクションを思い浮かべて行動を選択するか、再帰レベルによって見ていく。*B* は 1 つの運動に対し評価 (12) の高い 3 つの候補を比較することができる。Fig.7 に示すように、再帰レベルが 0 の時は、相手の運動に対して純粋に反応を返すだけである。一方、再帰レベルが 1 以上になると自分の反応に対する相手の反応、そのまた反応、… のようにインタラクションを予想し、現在の運動を選択するようになっている。

$n = 0$ の場合、*A* の初期動作「立ち上がる」(*su*) に対して *B* は「お辞儀をする」(*b*) を選択している。 $n = 1$ の場合も同様に *B* は「お辞儀をする」を選んでいるが、この時 *B* は自分がその動作を行った後に *A* が再び「立ち上がる」動作をするという予測を行っている。

$n = 2$ になると *B* は「手を叩く」という動作を選択する。これは、*B* が *A* の初期運動に対して「手を叩く」反応をした後に、*A* が「ハグ」し再び *B* が「手を叩く」というインタラクションを見出し、かつその評価が「立ち上がる」を選択したことで起こると考えるインタラクションより良いと判断したからである。

再帰計算によるインタラクションの予想では、行動選択が収束する場合・振動する場合に分けられる。収束する場合はさらに 2 通りに分けられ、図 7 のように評価の高い行動ループを発見して行動選択が収束する場合と、相手の初期運動に対して適切な解を見つけれず、*B* の選択が短期的な評価に偏って

いく場合である。前者の場合は先々のインタラクションに従って、後者の場合は目先の反応に従って選択される。*B* の選択が振動する場合は、再帰レベルが偶数の場合には最終的な自分自身のとりやすい運動の評価利得に従って、奇数の場合は相手のとりやすい運動の評価に従ってインタラクション決定される。

他者とのインタラクションを通じて他者を読んでいって自己の中での他者のモデルを確立し、それによって自己の振る舞いを決定する。人間が行っているプロセスを模倣していると考えられる。

4. おわりに

本稿での成果は以下の通りである。

1. HMM によって記号化された運動と自然言語を CRF によって連想させることによって、運動から言語、言語から運動を文脈に応じて柔軟に連想するモデルを構築した。CRF の素性関数には運動記号の変化と文節の変化を観測する関数を採用した。実験の結果、文脈が存在するところでは文脈に応じた連想を、そうでないところは単体の運動記号に基づく連想が行われた。
2. 提案した運動・言語相互連想モデルと IBM 翻訳モデル 2 を用いて、他者とのインタラクションを予想し、自己のとるべき運動を決定する手法を提案した。文長を評価に加えたところ、このモデルではインタラクションが長くなると選択される運動が収束する場合と振動する場合に分かれることが判明した。

参考文献

- [Brown 93] P. F. Brown, V. J. D. Pietra, S. A. D. Pietra, and R. L. Mercer: The Mathematics of Statistical Machine Translation: Parameter Estimation, Computational Linguistics, No. 2, pp. 263–311, Vol. 19 (1993).
- [Inamura 04] T. Inamura, I. Toshima, H. Tanie, and Y. Nakamura: Embodied Symbol Emergence Based on Mimesis Theory, The Intl. J. of Robotics Research, No. 4–5, pp. 363–377, Vol. 23 (2004).
- [伊藤 99] 伊藤昭: 心を読む能力の創発—マルチプレイヤー四人のジレンマゲーム—, Cognitive Studies, Vol. 6, No. 1, pp. 77–87 (1999).
- [Kudo 02] T. Kudo and Y. Matsumoto: Japanese Dependency Analysis using Cascaded Chunking, Proc. of the 6th Conf. on Natural Language Learning, pp. 1–7, Vol. 20 (2002).
- [Premack 08] D. Premack and G. Woodruff: Does the chimpanzee have a theory of mind?, The Behavioral and Brain Sciences, Vol. 1, No. 4, pp. 515–526 (1978).
- [高野 05] 高野雅典, 加藤正浩, 有田隆也: 心の理論における再帰レベルの進化に関する構成論的手法に基づく検討, 認知科学, No. 3, pp. 221–233, Vol. 12 (2005).
- [Takano 08] W. Takano and Y. Nakamura: Integrating Whole Body Motion Primitives and Natural Language for Humanoid Robots, IEEE-RAS Intl. Conf. on Humanoid Robots, pp. 708–713 (2008).