

参照点に依存した確率モデルの結合による動作の生成

○杉浦孔明 (京都大学, ATR), 岩橋直人 (情報通信研究機構, ATR)

Motion Generation by Combining Reference-Dependent Probability Models

*Komei Sugiura (Kyoto University, ATR), Naoto Iwahashi (NiCT, ATR)

Abstract— This paper proposes a method that learns motion concepts shown by users and generates a sequence of motion by combining reference-dependent probability models. This method estimates the reference point, the coordinate system, and the model parameters of the trajectory of the motion from observation data. Here the motion concepts are learned by using Hidden Markov Models (HMMs). Then, in motion generation phase, our method combines HMMs to generate trajectories to accomplish goal oriented tasks.

Key Words: human-robot interaction, HMM, motion generation, language acquisition

1. はじめに

日常的な環境において自然言語を用いた機械とのコミュニケーションを実現するためには、言語的・非言語的コンテキストから必要な特徴を抽出し、人間の用いるシンボルとマッピングする機構が必要である。このとき用いられるシンボルの指示対象は、ユーザによっても変わりうる。家庭で交わされる「食器棚からコップ出して」のような会話を例に挙げよう。機械にこのようなタスクを遂行させるためには、各家庭に固有の「食器棚」を認識させた上で、その「食器棚」に扉がついているか、取手があるかなどの形状により、適切な動作を選択させなければならない。このような背景から、動作とシンボル(動詞)、あるいはオブジェクトとシンボル(名詞や形容詞)を、ボトムアップに結びつける手法が研究されている [1, 2, 4].

動作のシンボル化に着目すると、このようなボトムアップ的手法は2つの重要性を持ち得る。第一に、ユーザから与えられたタスクを、ユーザが理解できる要素に分解し呈示できること、第二に、ユーザがシンボルを組合せて機械に指令を与えられること、である。特に前者は、安全性の観点から重要である。例えば、タスクがいくつかの動作からなるとき、機械が行なう予定の動作をユーザが理解しやすいシンボルを用いてあらかじめ説明すれば、ユーザはその動作を許可するかどうかを判断できる。

以上のような背景から、本論文では、参照点に依存した動作概念の学習および生成手法を提案する。具体的には、ロボットアーム、カメラ、マイクなどの構成要素からなるシステム (Fig. 1 [1]) を念頭におき、画像入力からの確率モデルの学習と、確率モデルを用いた軌道の生成手法について述べる。

2. 参照点に依存した空間的移動概念の学習と生成

2.1 参照点に依存した動作

空間的移動の概念には、参照点に依存しているものがある。例として、「ユーザが Fig. 2 に示す点線に沿って縫いぐるみを動かす」動作を考える。この動作は、参

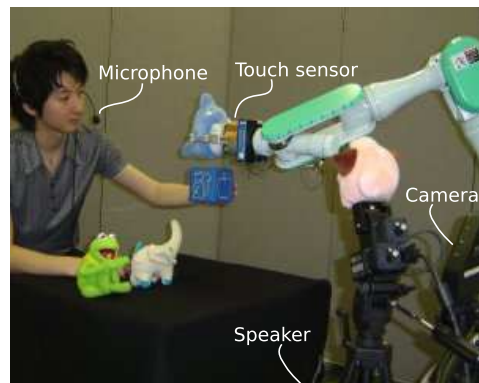


Fig.1 ヒューマンロボットインタラクション

照点を「画面右側の箱」とした場合には、「のせる」というラベルを与えることができる。すなわち、「ぬいぐるみを箱にのせた」という説明が妥当である。一方、参照点が「画面中央で静止している縫いぐるみ」である場合には、「とびこえさせる」というラベルが妥当である。

認知言語学では、外部世界を解釈する主体のプロセスにおいて焦点化される存在のうち、相対的に際立って認知される対象をトラジェクタ、これを背景的に位置付けるオブジェクトをランドマークとして区別する [6]. これにより、「～の左」や「～から離れた」など、対象の関係に基づく概念を記述している。

ここで、このような参照点に依存する動作の概念を、ユーザがラベルを発話した上で、物体を操作してロボットに学習させる問題を考える。このとき、適切な軌道をロボットに再現させるためには、(1) 参照点、(2) 参照点に依存した座標系、(3) その座標系での軌道生成パラメータ、の3つを推定しなければならない。例として、「とびこえさせる」と「ちかづける」の概念における、参照点と座標系について考えよう (Fig. 3). 図のように、「とびこえさせる」の概念は、参照点(ランドマーク)の位置を原点とする、カメラ座標系を平行移動した座標系を用いることができる。また、「ちかづける」の概念では、近づく対象をランドマークとし、ランドマークの位置を原点として、ランドマークからトラジェク

タへ向かう方向を x 軸とした直交座標系を選ぶのが妥当であろう。

谷らは、リカレントニューラルネットによる動作学習手法 [4] を提案しているが、参照点の推定を行っていないため、対象の位置が変化した場合に誤った軌道が生成される。また中村らは、参照点に依存しない動作を隠れマルコフモデル (HMM) を用いて生成させている [4]。我々は、参照点に依存する動作の概念を、軌道に関する運動情報 (座標, 速度, 加速度) の確率モデルとして学習する手法を提案している [1, 5]。以下では、本論文で提案する動作の学習および生成手法について説明する。

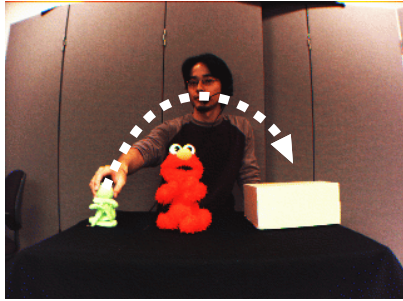


Fig.2 入力画像

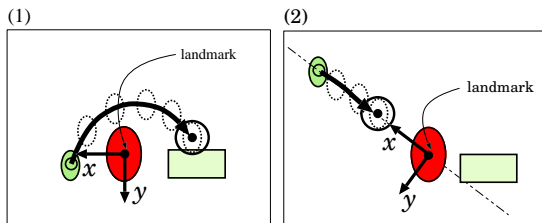


Fig.3 参照点に依存した座標系の表現: (1) 「とびこえさせる」、(2) 「ちかづける」

2.2 確率モデルを用いた動作の学習

移動するオブジェクトが一つで複数の静止オブジェクトが存在する動画画像から抽出された特徴量データ集合 $V = \{V_1, V_2, \dots, V_N\}$ が学習データとして与えられるものとする。ただし、各動画画像 V_i は、2次元カメラ座標系における、トラジェクタの運動情報 \dagger_i および静止オブジェクトの位置の集合 O_i からなるとする。運動情報 (軌道に関する情報) としては、トラジェクタの位置 \mathbf{x}_t ・速度 $\dot{\mathbf{x}}_t$ ・加速度 $\ddot{\mathbf{x}}_t$ の時系列を用いる。すなわち、 \mathcal{Y}_i は、 $\mathbf{y}_t = [\mathbf{x}_t, \dot{\mathbf{x}}_t, \ddot{\mathbf{x}}_t]^T$ の時系列である。また、各静止オブジェクトと軌道開始点と画像中心点を参照点の候補とし、その集合を $L_i = \{l_1^i, l_2^i, \dots, l_{M_i}^i\}$ とする。

各 V_i に設定する 2次元座標系は、参照点を原点とし、参照点および軌道開始点の情報をもとに設定されるものとする。ただし、座標軸の向き k の候補は K 種類であるとする。座標軸の向きは概念に固有であるが、参照点は各学習データ V_i に対して選択される。

ここで、 \mathcal{Y}, k と参照点 l によって決定される座標系でのトラジェクタ軌道を、 $F(\mathcal{Y}, k, l)$ と表すことにする。この設定のもと、最適な座標軸の設定方法 \hat{k} および参照点を探索しながら、軌道に関する確率モデルのパラ

メータ λ を尤度最大化基準により学習する。つまり、

$$(\hat{\lambda}, \hat{k}, \hat{m}) = \operatorname{argmax}_{\lambda, k, m} \sum_{i=1}^N \log P(F(\mathcal{Y}_i, k, l_{m_i}^i); \lambda). \quad (1)$$

ここで、 m_i は参照点 $l_{m_i}^i$ を選択することを示す。また、 $\mathbf{m} = (m_1, m_2, \dots, m_N)$ とする。確率モデルとして隠れマルコフモデル (HMM) を用いた場合の解法の詳細については、[5] を参照されたい。

2.3 確率モデルの結合による動作の生成

本節では、確率モデルを結合して軌道を生成する方法を提案する。提案手法では、2種類の生成を扱う。すなわち、(1) 「トラジェクタ初期位置 \mathbf{x}_0 と目標位置 \mathbf{x}_n 」を入力として、最も尤度の高い「<動作名, ランドマーク>の列 \hat{A} , および軌道 $\hat{\mathcal{Y}}$ 」を生成する、(2) トラジェクタおよび「<動作名, ランドマーク>の列 A 」を入力として、最も尤度の高い軌道 $\hat{\mathcal{Y}}$ を生成する、の2通りである。

ここで、確率モデルとして HMM を用いて、軌道の生成を行なうことを考える。HMM による音声合成のように、基準となる座標系を共有した上で2つの HMM を結合する手法は広く行なわれている。しかし、前節の手法を用いて得られた確率モデルは、それぞれの座標軸上での軌道を表しているの、そのまま結合することはできない。例えば、Fig. 4 のように「あげる」と「ちかづける」を表す HMM を結合する場合は、「あげる」の終点をトラジェクタ位置として、「ちかづける」を表す HMM のパラメータ λ を座標変換する必要がある。

いま、 A 中に含まれる動作名称列が、 n 個の動作パラメータ列 $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)$ を表すとする。このとき λ_j を以下のように変換する。

$$E_{\mathbf{y}_s}(\lambda_j) = \mathcal{W}_{k,l} \left(\begin{bmatrix} E_{\mathbf{x}_s}(\lambda_j) - E_{\mathbf{x}_0}(\lambda_j) \\ E_{\dot{\mathbf{x}}_s}(\lambda_j) \\ E_{\ddot{\mathbf{x}}_s}(\lambda_j) \end{bmatrix} \right) + E_j' \quad (2)$$

$$E_j' = [E_{\mathbf{x}_{S_{j-1}}}(\lambda_{j-1}), \mathbf{0}, \mathbf{0}]^T \quad (3)$$

$$E_{\mathbf{x}_{S_{j-1}}}(\lambda_0) = \mathbf{x}_0 \quad (4)$$

$$E_{\mathbf{x}_{S_n}}(\lambda_n) = \mathbf{x}_n \quad (5)$$

ただし、 s は HMM の状態を表し、 S_j は λ_j の最終状態であるとする。また、 $E_{\mathbf{x}_s}(\lambda_j)$ は、 λ_j のうち状態 s における平均特徴ベクトルのうち位置に関するものを表し、 $\mathcal{W}_{k,l}$ は、参照点を l とする座標系 k から世界座標系への変換を表す。

次に、分散に関するパラメータ $V_{\mathbf{y}_s}(\lambda_j)$ を以下に従い変換する。

$$V_{\mathbf{y}_s}(\lambda_j) = V_{\mathbf{y}_s}(\lambda_j) + V_j' \quad (6)$$

$$V_j' = \begin{cases} \mathbf{0} & \text{if } l \in O_i, \\ [V_{\mathbf{x}_{S_{j-1}}}(\lambda_{j-1}), \mathbf{0}, \mathbf{0}]^T & \text{otherwise} \end{cases} \quad (7)$$

つまり、動作概念がランドマークに依存しない場合、動作系列の分散パラメータ $V_{\mathbf{y}_s}(\lambda_j)$ は直前の動作概念の最終状態における分散パラメータ $V_{\mathbf{y}_{S_{j-1}}}(\lambda_{j-1})$ に依存

する。これにより、ランドマークに依存した動作を確実に達成しつつ、軌道を滑らかに結合することができる。ただし上式では、分散の回転を考えていない。これは、HMM 合成では共分散行列の対角成分のみを用いるのが一般的なためであるが、厳密には他の成分も考える必要がある。

このようにして変換された λ_j を結合したものを Λ' とする。よって上記の (1) のタイプの生成においては、軌道開始点 x_0 および終了点 x_n の制約条件のもとでの、最尤軌道 \hat{y} の推定を行なう。すなわち、深さ n までの範囲において以下の探索を行なう。

$$(\hat{y}, \hat{\Lambda}') = \underset{y, \Lambda'}{\operatorname{argmax}} \log P(y; \Lambda') \quad (8)$$

上式の解は、[3] で提案されている最適化法によって求めることができる。

なお、(2) のタイプの生成は上の流れと同様であるが、 x_n を設定しない。さらに、 A が与えられるので、 Λ' の探索部分を省略する。

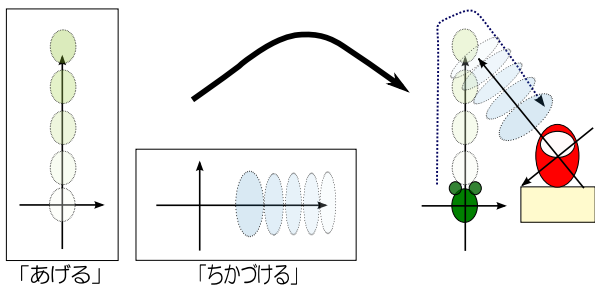


Fig.4 HMM の結合における座標変換

3. 実験

3.1 実験設定

まず動作生成に用いる動作要素を用意する。本実験では、動作要素をユーザの教示により学習させる。そのために、以下の7個の概念を学習させた。

「あげる」、「ちかづける」、「はなす」、「まわす」、「のせる」、「さげる」、「とびこえさせる」

座標系は以下の k_1 から k_4 のいずれかとする。

- k_1 : ランドマークを原点とする、カメラ座標系を平行移動した座標系。ただし、変換後の座標系において軌道開始点の x 座標が負になる場合には、さらに x 軸を反転させる。
- k_2 : ランドマークを原点とし、軌道開始点に向かう軸を x 軸とする直交座標系。
- k_3 : 軌道開始点を原点とする、カメラ座標系を平行移動した座標系。
- k_4 : 画面中心を原点とする、カメラ座標系を平行移動した座標系。

Fig. 5 に学習データの例と認識結果 (参照点, 座標系) を示す。

(1) のタイプの軌道の生成においては、探索の深さ $n = 3$ とする。つまり、生成される連結動作は最大 n 個の動作要素からなる。また、軌道の生成においてオブジェクトとの衝突は考えない。

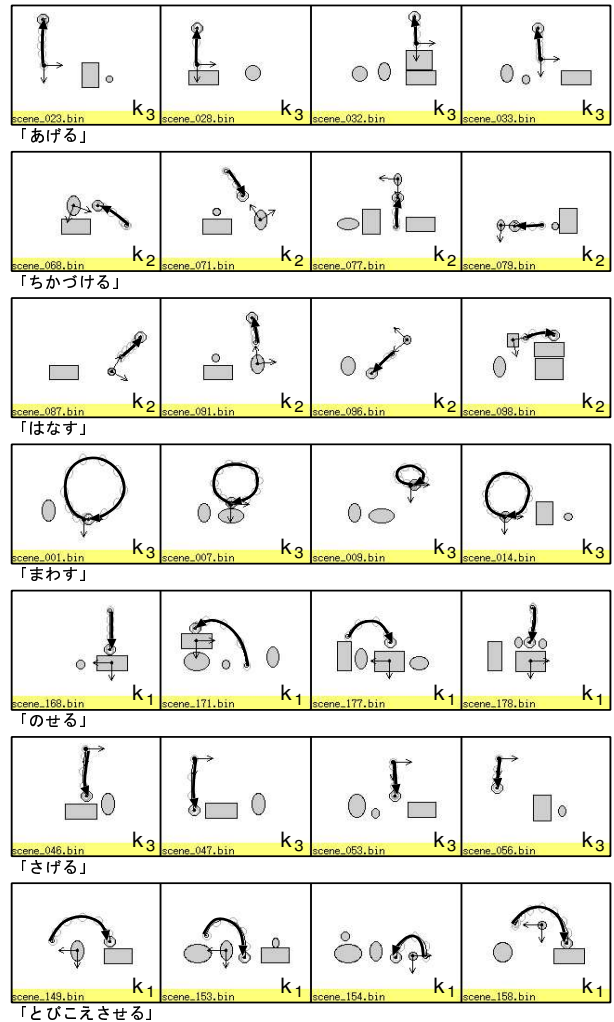


Fig.5 学習データの例と学習の結果選択された座標系

3.2 軌道生成

3.2.1 終点を指定した軌道生成

Fig. 6 に示すような5個のオブジェクトが存在する環境において、オブジェクト1をトラジェクタとしてカメラ座標系の格子点64個に対して軌道を生成させた。図には生成された軌道の例を示した。図において、 $\langle \rangle$ は最も尤度の高い動作名とランドマークの組を示す。例えば、 \langle のせる, 2 \rangle \langle ちかづける, 5 \rangle は、終点までの動作が2つの動作の組合せとして生成されたことを示し、最初の動作は「1を2にのせる」こと、次の動作は「1を5にのせる」ことを意味している。

図より、 \langle のせる, 2 \rangle \langle ちかづける, 5 \rangle と \langle のせる, 2 \rangle \langle はなす, 2 \rangle は途中まで同じ軌道をとることがわかる。これは、「のせる」を学習した確率モデルにおいて、最終位置に対する分散が小さいことが学習されたため起こったと考えられる。これに対し、 \langle はなす, 2 \rangle と \langle はなす, 2 \rangle \langle とびこえさせる, 4 \rangle の軌道のように、「はなす」動作を共有していても、実際の軌道が異なる確率モデルも存在する。これは、「はなす」として学習された確率モデルの最終位置の分散が、「のせる」の場合よりも大きいためであると考えられる。

本手法により64通りの目標点に全て到達することが

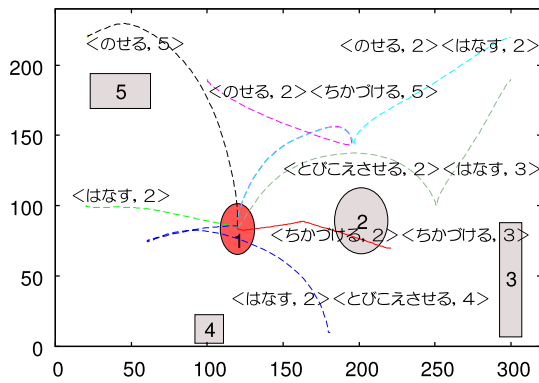


Fig.6 終点を指定した軌道生成の例

できたが、本手法では衝突判定を行っていないので、<くちかづける, 2><くちかづける, 3>のように、物体上を通過する軌道が生成され得る (Fig. 6 参照)。これを避けるためには、得られた最尤軌道を変化させる手法、衝突しない軌道のうち最も尤度が高い軌道を用いる手法などが考えられる。

3.2.2 動作列を指定した軌道生成

次に、動作列を指定した場合の軌道生成の実験結果について述べる。Fig. 7 は 2 種類の動作指令、「(a) オブジェクト 1 を、オブジェクト 2 の上を飛び越えさせて、下げて、オブジェクト 4 に近づける」動作と「(b) オブジェクト 2 を、オブジェクト 1 の上を飛び越えさせた後、再びオブジェクト 1 の上を飛び越えさせ、オブジェクト 5 に載せる」動作を示したものである。このとき動作列は、

- (a) トラジェクタ = 1, $A = \langle \text{とびこえさせる, 2} \rangle \langle \text{さげる, ランドマークなし} \rangle \langle \text{くちかづける, 4} \rangle$
- (b) トラジェクタ = 2, $A = \langle \text{とびこえさせる, 1} \rangle \langle \text{とびこえさせる, 1} \rangle \langle \text{のせる, 5} \rangle$

とした。図より、3つの動作を滑らかに結合できていることがわかる。さらに、この結果を定量的に検討するために、Fig. 8 に、(a) の場合の位置・速度・加速度の変化を示す。この結果も、提案手法が動作を滑らかに結合できる、という結論を支持している。

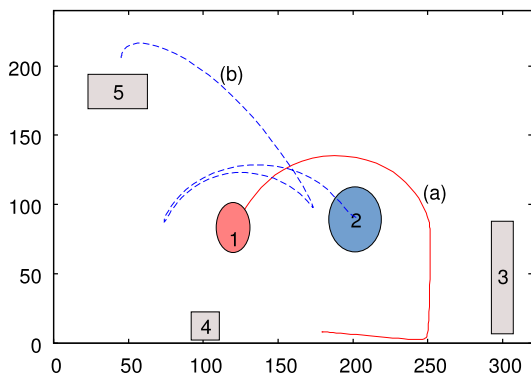


Fig.7 動作列を指定した軌道生成の例

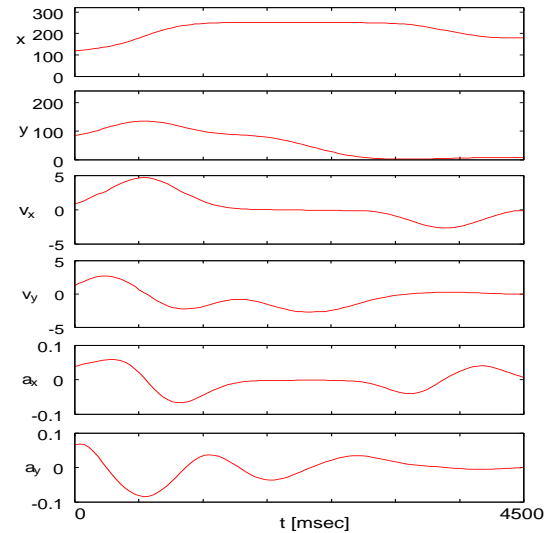


Fig.8 Fig. 7(a) の軌道における座標, 速度, 加速度の変化

4. おわりに

人間と機械(ロボット)が共存する環境において、機械が内部状態を人間にわかりやすく伝える機能は重要であり、ヒューマンインタフェースの分野などでも研究が行なわれている。例えば、「物を運搬する移動ロボットが次にどのような動作をするかを周囲の人間に伝える」機能は安全の観点から有意義である。本研究では、ユーザの行為を実世界にグラウンドした知識として獲得し、学習した行動要素群を用いて目標への移動を達成する軌道生成法について述べた。応用例としては、オフィス環境において、ユーザにより指定された目標位置までの経路を生成し、ユーザに動作系列を自然言語で伝えるロボットなどが挙げられる。

謝辞

本研究は、日本学術振興会科学研究費補助金 18・2972 および国立情報学研究所共同研究「能動的ハンドインタラクションによる実世界言語コミュニケーションの学習に関する研究」による研究助成を受け実施したものである。

参考文献

- [1] Iwahashi, N.: Robots That Learn Language: Developmental Approach to Human-Machine Conversations, *Symbol Grounding and Beyond: Proceedings of the Third International Workshop on the Emergence and Evolution of Linguistic Communication* (Vogt, P. et al.(eds.)), Springer, pp. 143-167 (2006).
- [2] Roy, D.: Grounding Words in Perception and Action: computational insights, *Trends in Cognitive Science*, Vol. 9, No. 8, pp. 389-396 (2005).
- [3] Tokuda, K., Kobayashi, T. and Imai, S.: Speech parameter generation from HMM using dynamic features, *Proceedings of International Conference on Acoustics, Speech, and Signal Processing*, pp. 660-663 (1995).
- [4] 土井利忠, 藤田雅博, 下村秀樹 (編): 脳・身体性・ロボット, シュプリンガー・フェアラーク東京 (2005).
- [5] 羽岡哲郎, 岩橋直人: 言語獲得のための参照点に依存した空間的移動の概念の学習, 信学技報, PRMU2000-105, pp. 39-46 (2000).
- [6] 山梨正明: 認知言語学原理, くろしお出版 (2000).