

距離・画像センサと複素相関を用いた 回転不変なジェスチャ認識

小林 宏章¹ 加藤 秀章¹ 尺長 健¹

概要: 距離・画像センサを用いたジェスチャ認識について述べる。従来、我々は距離・画像センサデータを用いた人物姿勢追跡法を提案している。本稿では、この方法によって時々刻々得られる人物上半身の姿勢を、重力方向まわりの回転に不変な類似度を複素相関を用いて構成することにより、胴体の回転に不変なジェスチャ認識法を提案する。また、実データに適用した結果を報告する。

1. はじめに

ジェスチャ認識は、人間と計算機間の自然なユーザインタフェースとして注目されている [1][2]。映像に基づくジェスチャ認識においては、従来、背景と人物の分離が難しかったことから、顕著な画像特徴に重点をおいた追跡 [3][4] を試みたり、動作に関する学習を利用したトップダウン型の認識系が提案されてきた。一方、距離・画像センサが容易に利用できるようになったのに伴い、距離情報を併用することで、容易に人物領域を抽出でき、さらに距離情報を用いることで人物の追跡の精度も向上し、より正確なジェスチャ認識が行えると考えられる。四宮ら [5] は距離・画像センサを用いた追跡系を構成し、人物の上半身を 14 自由度の姿勢空間で取扱うことで 3 次元疎テンプレート追跡に基づく姿勢追跡を提案している。本稿では、この姿勢追跡により得られる時系列姿勢情報を用いてジェスチャ認識系を構成する。このために、まず、姿勢間の類似度、および、ジェスチャ間の類似度を定義する。本稿では、重力方向（通常姿勢では、胴体の軸とほぼ一致する）まわりの回転に不変な姿勢類似度を複素相関を用いて定義し、これをベースとして、胴体の回転に不変なジェスチャ認識を目指す。第 2 章では、上半身の関節物体としての取扱いと、3 次元姿勢追跡系の概要を述べる。第 3 章では、まず、複素相関を用いて姿勢類似度、および、ジェスチャ類似度を定義する。また、姿勢追跡により得られたパラメータからジェスチャを認識する方法について述べる。第 4 章では本稿で提案するジェスチャ認識を用いた実動画実験の結果を述べる。カメラに対して正面を向いてジェスチャを行う動

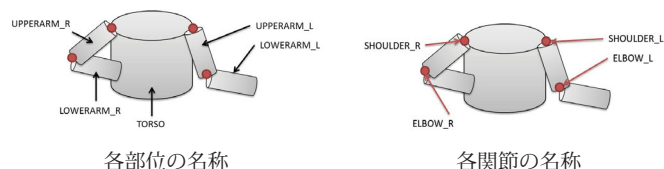


図 1 関節物体モデル

画像だけではなく、15, 30, 45 度の各方向を向いてジェスチャを行う動画像においても実験を行い、結果を述べる。

2. 距離・画像センサによる関節物体の姿勢追跡

2.1 関節物体の姿勢

本稿で追跡の対象とするのは、人間の上半身である。本稿では、上半身は胴体、左上腕、左下腕、右上腕、右下腕の 5 つの剛体が左肩、左肘、右肩、右肘の 4 つの関節によって連結されているものとする。この上半身モデルの各部位と各関節を図 1 示す。なお、左肩、右肩の回転自由度を 2 とし、左肘、右肘の回転自由度もまた 2 とする。

2.2 座標系の定義

一般に、関節物体の姿勢は、主となる物体の姿勢と、それぞれの関節位置、関節角度により定義される。

本稿では、上半身モデルのレンダリングや姿勢推定問題を考えるため、以下に示す 4 つの座標系を考える。

- (1) M_C : カメラ (中心) 座標系
- (2) M_O : 胴体 (中心) 座標系
- (3) M_U : 上腕 (中心) 座標系
- (4) M_L : 下腕 (中心) 座標系

なお、 M_U と M_L については左右に別々に座標系を考える必要があるため、実装は 6 個の座標系を考える必要がある。本稿では、左右の相互作用は考慮しないため、以下で

¹ 岡山大学
Okayama University

は4つの座標系を用いた議論を行う。

2.3 3次元サーフィスモデル

本稿では、追跡対象である3次元物体はサーフィスモデルで取り扱う。サーフィスモデルでは、物体の表面を小さなパッチの集合として扱う。ここで、物体モデルの表面を構成する点を胴体座標系、上腕座標系、下腕座標系を用いて表す。

一方、画像生成過程はカメラ座標で記述される。ある姿勢の物体モデルから画像を生成するためには、胴体/上腕/下腕座標からカメラ座標 $M_C = [X_C Y_C Z_C]^T$ 変換する必要がある。表記の簡単化のために、 M_C, M_O, M_U, M_L の同次座標表現 $\tilde{M}_C = [M_C^T 1]^T, \tilde{M}_O = [M_O^T 1]^T, \tilde{M}_U = [M_U^T 1]^T, \tilde{M}_L = [M_L^T 1]^T$ を用いる。

2.3.1 胴体の表現と姿勢行列

胴体については、6次元姿勢(3次元並進+3次元回転)を 4×4 姿勢行列で取り扱う。胴体座標系で表される3次元点の同次表現を $\tilde{M}_O = [M_O^T 1]^T$ とし、胴体座標系からカメラ座標系への変換 D_{OC} を胴体の姿勢行列と考えることができる。

$$D_{OC} = \begin{bmatrix} R_{OC} & t_{OC} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (1)$$

$$\tilde{M}_C = D_{OC} \tilde{M}_O \quad (2)$$

ここで式(2)における R_{OC} は胴体の回転を表す行列であり、 t_{OC} は胴体の並進を表すベクトルである。

2.3.2 上腕の表現と姿勢行列

上腕の姿勢は、胴体の姿勢 D_{OC} と肩を中心とした回転の合成により得られる。上腕座標系で表される3次元点の同次表現を M_U とし、上腕座標系から胴体座標系への変換を D_{UO} と表記する。胴体座標系における肩関節の位置を t_S とする。ここで、肩関節は上腕座標系の原点にあるものとする。また、肩での回転を表す回転行列を R_S とすると、上腕座標で表される点のカメラ座標系への変換は以下の式(4)で表せる。

$$D_{UO} = \begin{bmatrix} R_S & t_S \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (3)$$

$$\tilde{M}_C = D_{UC} \tilde{M}_U \quad (4)$$

ここで、

$$D_{UC} = D_{OC} D_{UO} \quad (5)$$

$$R_S = \begin{bmatrix} c_{\theta_S} & 0 & s_{\theta_S} \\ 0 & 1 & 0 \\ -s_{\theta_S} & 0 & c_{\theta_S} \end{bmatrix} \begin{bmatrix} c_{\phi_S} & -s_{\phi_S} & 0 \\ s_{\phi_S} & c_{\phi_S} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (6)$$

なお、 $c_{\theta_S} = \cos \theta_S, s_{\theta_S} = \sin \theta_S$ とする。

2.3.3 下腕の表現と姿勢行列

下腕の姿勢は、上腕の姿勢 D_{UC} と、肘を中心とした回

転の合成により得られる。下腕座標系で表される3次元点の同次表現を M_L とし、下腕座標系から上腕座標系への変換を D_{LU} と表記する。上腕座標系における肘の位置を t_E とする。ここで、肘関節は下腕座標系の原点にあるものとする。さらに肘での回転(2自由度)を表す回転行列を R_E とすると、下腕の点をカメラ座標系へ変換する式(8)は以下のように表せる。

$$D_{LU} = \begin{bmatrix} R_E & t_E \\ \mathbf{0}^T & 1 \end{bmatrix} \quad (7)$$

$$\tilde{M}_C = D_{LC} \tilde{M}_L \quad (8)$$

ここで、

$$D_{LC} = D_{UC} D_{LU} \quad (9)$$

$$R_E = \begin{bmatrix} c_{\phi_E} & -s_{\phi_E} & 0 \\ s_{\phi_E} & c_{\phi_E} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} c_{\theta_E} & 0 & s_{\theta_E} \\ 0 & 1 & 0 \\ -s_{\phi_E} & 0 & c_{\theta_E} \end{bmatrix} \quad (10)$$

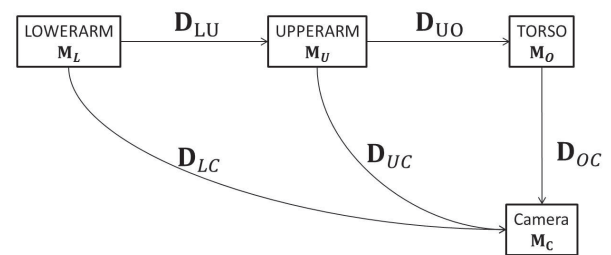


図2 変換の関係

なお、 $c_{\theta_E} = \cos \theta_E, s_{\theta_E} = \sin \theta_E$ とする。4つの座標系の相互変換をまとめて図2に示す。座標系 M_L, M_U, M_O からカメラ座標系 M_C への座標変換は、それぞれ D_{LC}, D_{UC}, D_{OC} で行える。また、図中の座標変換の逆変換は逆行列で表される。例えば、 $D_{CL} = D_{LC}^{-1}$ である。

2.4 基準姿勢における3次元疎テンプレートの作成

既知形状物体の3次元姿勢追跡は直前の物体の見えから作成された疎テンプレートによる追跡に基づいて行う。以下に、その詳細を述べる。前時刻 $t-1$ の推定姿勢が与えられた場合、この姿勢から姿勢行列 D_{t-1} を算出できる。以下では、前時刻の推定姿勢から予想される姿勢を基準姿勢とよぶ。基準姿勢から生成される姿勢行列を基準姿勢行列と呼び、 \hat{D} と表す。(簡単には、 $\hat{D} = D_t$ としてよい。) 通常の Z-buffer 法の拡張として、画像 $I(u, v)$ 、距離 $Z(u, v)$ を更新する際に、同時に $X(u, v), Y(u, v)$ の値も更新することにより各画素 (u, v) において画素値 $I(u, v)$ に対応する3次元座標 $[X(u, v) Y(u, v) Z(u, v)]$ を計算する。これにより、基準姿勢行列 \hat{D} によって作成されるテンプレート画像 $\hat{I}(u, v)$ から3次元疎テンプレートを作成できる。すなわち、局所領域内で輝度値が最大・最小となる点を選択す

ることにより、疎テンプレート $\{(u_j, v_j)\}$ を作成し、この各点 (u_j, v_j) に対応する 3 次元座標を M_j とすると、3 次元疎テンプレートは次のように記述できる。

$$\{(M_j; \hat{I}_j)\} = \{(X_j Y_j Z_j; \hat{I}_j)\} \quad (11)$$

ただし、 $X_j = X(u_j, v_j)$, $Y_j = Y(u_j, v_j)$, $Z_j = Z(u_j, v_j)$, $\hat{I}_j = \hat{I}(u_j, v_j)$ である。

2.5 3次元疎テンプレートマッチング

3次元疎テンプレートを構成する j 番目の点 M_j が胴体、上腕、下腕のどこに属するかによって、式 (12) によって画像上の位置 m_{jk} を算出することになる。すなわち、各部位において括弧内を定数として取り扱うことができるため、 k 番目のサンプル姿勢に対応する姿勢変換行列 $\delta D_k, \delta D'_k, \delta D''_k$ を用いて、画像座標 m_{jk} が計算される。

$$\begin{cases} m_{jk} = (A)\delta D_k(D_{OC}\tilde{M}_j) \\ m_{jk} = (AD_{OC})\delta D'_k(D_{UO}\tilde{M}_j) \\ m_{jk} = (AD_{UC})\delta D''_k(D_{LU}\tilde{M}_j) \end{cases} \quad (12)$$

同様に、距離 Z_{jk} は次式によって求められる。

$$\begin{cases} (X_{jk}Y_{jk}Z_{jk}1)^\top = \delta D_k(D_{OC}\tilde{M}_j) \\ (X_{jk}Y_{jk}Z_{jk}1)^\top = (D_{OC})\delta D'_k(D_{UO}\tilde{M}_j) \\ (X_{jk}Y_{jk}Z_{jk}1)^\top = (D_{UC})\delta D''_k(D_{LU}\tilde{M}_j) \end{cases} \quad (13)$$

ただし、 \tilde{M}_j は M_j の同次座標を示す。また、 $\delta D_k, \delta D'_k, \delta D''_k$ はサンプル姿勢の 14 パラメータ $[\alpha_1 \alpha_2 \dots \alpha_{14}]^\top = [\delta\theta_0 \delta\phi_0 \delta\psi_0 \delta t_1 \delta t_2 \delta t_3 \delta\phi_S \delta\theta_S \delta\phi_E \delta\theta_E \delta\phi'_S \delta\theta'_S \delta\phi'_E \delta\theta'_E]$ を用いて次式で示される。

$$\delta D_k = \begin{bmatrix} c_2c_3 & -c_2s_3 & s_2 & \delta t_1 \\ s_1s_2c_3 + c_1s_3 & c_1c_3 - s_1s_2s_3 & -s_1c_2 & \delta t_2 \\ s_1s_3 - c_1s_2c_3 & c_1s_2s_3 + s_1c_3 & c_1c_2 & \delta t_3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (14)$$

$$\delta D'_k = \begin{bmatrix} c_jc_{j+1} & -c_js_{j+1} & s_j & 0 \\ s_{j+1} & c_{j+1} & 0 & 0 \\ s_js_{j+1} & s_js_{j+1} & c_j & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (15)$$

$$\delta D''_k = \begin{bmatrix} c_{j+2}c_{j+3} & -s_{j+3} & s_{j+3}c_{j+3} & 0 \\ c_{j+2}s_{j+3} & c_{j+3} & s_{j+2}s_{j+3} & 0 \\ -s_{j+2} & 0 & c_{j+2} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (16)$$

ここで、 $c_i = \cos\alpha_i$, $s_i = \sin\alpha_i$ を示す。式 (15), (16) 中において、右手系では $j = 7$, 左手系では $j = 11$ である。

2.6 粗密追跡法

追跡では、姿勢サンプル集合の一部をフレームからフレームへと伝搬させる。ガウスノイズが各サンプルのすべての姿勢パラメータに付加され、その生成されたサンプルから

姿勢行列を算出し、最終的にすべてのサンプルは疎テンプレートマッチングによって評価される [6]。本稿では、粗密追跡法 [7] を用いることで、サンプル数の削減を図る。粗密追跡法では、ガウスノイズの標準偏差を段階的に減らすことで、効率よく姿勢空間を探索する。

2.7 カスケード姿勢推定

本稿では、人体 (上半身) の姿勢は胴体に 6 自由度、肩に 2 自由度、肘に 2 自由度をもつものとしている。したがって上半身の追跡には 14 次元の姿勢空間を取り扱う必要がある。

本稿では、胴体、右腕、左腕の順に姿勢を推定することで、姿勢空間を小さく抑えながら高次元の姿勢空間を探索する。まず、胴体の姿勢の推定を行い、胴体姿勢が推定された後、右腕の姿勢 (4 自由度) ついて姿勢サンプルを生成し、姿勢を推定する。最後に、胴体と右腕の姿勢を用いて左腕のサンプル姿勢を 4 次元姿勢空間で作成し、姿勢を推定する。各段の計算では、基準姿勢のまわりでサンプル姿勢に対応する疎テンプレートの位置を効率よく計算する。

2.8 距離・画像センサを用いたテンプレートマッチング

本稿の追跡では、3次元テンプレートを用いるため、画像情報と同時に距離情報を利用できる。これを利用し、距離・画像センサから得られた距離情報と物体モデルの各サンプルにおける奥行きを比較することで、画像情報のみを用いた追跡よりも高性能な追跡を期待できる。

2.8.1 距離情報を用いた人物追跡と背景領域の分離

画像情報では人物領域と背景領域の分離を厳密に行うことは一般には容易でない。しかし、距離情報が追跡時に利用できる場合には、距離のギャップを利用して簡単に両者を分離できる。すなわち人物領域と背景領域では、明らかに距離情報にギャップがあり、また、背景の方がカメラから遠いことを利用し、人物領域のみを抽出し、マスクをかけることができる。これにより見え画像によるテンプレートマッチングの精度が向上すると考えられる。

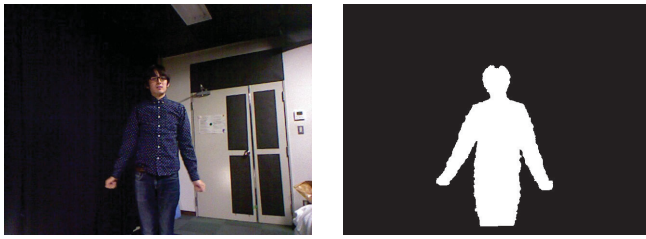
ここでは $d_0 \sim d_1$ [mm] の範囲内で人物が観測されるとする。各画素 $I(x, y)$ に対応する距離を $d(x, y)$ とすると、人物領域のマスク $m(x, y)$ は以下の式 (17) により表される。

$$m(x, y) = \begin{cases} 1 & d_0 \leq d(x, y) \leq d_1 \\ 0 & otherwise \end{cases} \quad (17)$$

なお、実装では $d_0 = 500$, $d_1 = 2400$ としている。見え画像とこの処理により抽出した人物領域の一例を図 3 に示す。

2.8.2 マスクを用いたテンプレートマッチング

前節で作成したマスク $m(u, v)$ を用いて、入力画像 I でテンプレートマッチングを行い、画像情報の評価値 ϵ_k を求める。



見え画像 人物マスク領域
図 3 見え画像と人物マスク領域

$$\epsilon_k = \frac{1}{\sum_{j=1}^{j_{max}} \rho(\beta I(u_{jk}, v_{jk})m(u_{jk}, v_{jk}) - \alpha \hat{I}_j)} \quad (18)$$

ただし、

$$\alpha = \frac{1}{\sum_{j=1}^{j_{max}} \hat{I}_j}, \quad \beta = \frac{\alpha \sum_{j=1}^{j_{max}} \hat{I}_j^2}{\sum_{j=1}^{j_{max}} \hat{I}_j I(u_{jk}, v_{jk})}$$

$$\rho(x) = \frac{x^2}{c^2 + x^2}$$

なお、 α はテンプレート画像の正規化係数、 β は入力画像の正規化係数、 $\rho(x)$ は Geman-McClure 関数である。また c は定数であり、ここでは、 $c = 0.5/n$ とする。(ただし n はテンプレートの点数)

2.8.3 距離情報を利用したテンプレートマッチング

サンプル姿勢 $\delta \mathbf{D}_k$ が与えられたとき、入力距離情報の評価値 ϵ'_k は、テンプレートの j 番め点の距離情報 Z_{jk} と入力距離情報 $d(j)$ を用いて次式で計算される。

$$\epsilon'_k = \frac{1}{\sum_{j=1}^{j_{max}} \rho'(Z'_{jk}m(u_{jk}, v_{jk}) - d(j))} \quad (19)$$

ただし、

$$\rho'(x) = \frac{x^2}{c'^2 + x^2}$$

$\rho'(x)$ は Geman-McClure 関数である。ここで、 c' は定数であり、適切な値を設定することでノイズの影響を低減することができる。距離・画像センサの出力が $1mm$ を単位とする整数値であることを勘案し、予備実験により、 $c' = 5.0$ とした。

2.8.4 画像情報と距離情報を併用した評価値の算出

画像情報に対して算出した評価値 ϵ_k と距離情報に対して算出した評価値 ϵ'_k を用いて姿勢の推定を行うことで、より適切な姿勢を決定できると考えられる。

あるサンプル姿勢における評価値 ϵ を両者の重み付き和によって次式で定義する。

$$\epsilon = (1 - \lambda)\epsilon_k + \lambda\epsilon'_k \quad (20)$$

3. 姿勢追跡を用いたジェスチャ認識

本章では、姿勢情報からジェスチャ認識を行う手法について述べる。各ジェスチャをあらかじめ複数回撮影し、各追跡結果として、得られる姿勢の時系列平均を登録ジェス

チャとする。これにより得られる登録ジェスチャと入力における姿勢情報とを評価し、入力動画中の人物が登録されているジェスチャのうち、どのジェスチャを行っているかを認識する。

3.1 ジェスチャ認識の流れ

まず、姿勢追跡によって得られた現在の姿勢と、各登録ジェスチャの開始フレームの平均姿勢とを照合し、ジェスチャの開始を検出する。次に、同様の方法でジェスチャ終了の検出を行う。開始と終了が検出された時点で、登録ジェスチャと入力ジェスチャの対応付けを行う。これは、人物がジェスチャを行うときロバストに認識するためである。最後に、対応付けにより得られた入力ジェスチャと登録ジェスチャとを照合し、そのジェスチャであるかどうかを判定する。

3.2 登録ジェスチャ

ジェスチャの登録には、胴体を除いた両腕の姿勢(肩、肘を結んだ3次元ベクトルと肘、手を結んだ3次元ベクトル)を用いる。第 g 登録ジェスチャを \mathbf{S}_g と表す。 \mathbf{S}_g は 30 フレームの姿勢(12次元)で表される 360 次元である。また、登録ジェスチャの数を G とする。また、第 g 登録ジェスチャの第 $f(f = 1, 2, \dots, F)$ フレームを \mathbf{S}_{gf} で表す。本稿では、登録ジェスチャのフレーム $F = 30$ とする。

以下に、登録ジェスチャの要素として扱うベクトルについて説明する。なお、 \mathbf{D}_{UC} は上腕座標からカメラ座標への変換行列であり、 \mathbf{D}_{OC} は胴体座標からカメラ座標への変換行列である。詳しくは論文 [5] で定義されている。上腕座標における肘の座標 $\widetilde{\mathbf{M}}_{UE}$ を以下の式でカメラ座標 $\widetilde{\mathbf{M}}_{CE}$ にする。

$$\widetilde{\mathbf{M}}_{CE} = \mathbf{D}_{UC}\widetilde{\mathbf{M}}_{UE} \quad (21)$$

同様に、胴体座標における肩の座標 $\widetilde{\mathbf{M}}_{OS}$ をカメラ座標 $\widetilde{\mathbf{M}}_{CS}$ にする。

$$\widetilde{\mathbf{M}}_{CS} = \mathbf{D}_{OC}\widetilde{\mathbf{M}}_{OS} \quad (22)$$

\mathbf{M}_{CS} と \mathbf{M}_{CE} より上腕の3次元ベクトル \mathbf{V}_1 を以下の式で求める。

$$\mathbf{V}_1 = \frac{\mathbf{M}_{CE} - \mathbf{M}_{CS}}{|\mathbf{M}_{CE} - \mathbf{M}_{CS}|} \quad (23)$$

下腕座標における手の座標 $\widetilde{\mathbf{M}}_{UH}$ を以下の式で上腕座標 $\widetilde{\mathbf{M}}_{CH}$ にする。

$$\widetilde{\mathbf{M}}_{CH} = \mathbf{D}_{UC}\widetilde{\mathbf{M}}_{UH} \quad (24)$$

式 (21) で求めた肘の座標 \mathbf{M}_{CE} と式 (24) で求めた手の座標 \mathbf{M}_{CH} より、下腕の3次元ベクトル \mathbf{V}_2 を以下の式で求める。

$$\mathbf{V}_2 = \frac{\mathbf{M}_{CH} - \mathbf{M}_{CE}}{|\mathbf{M}_{CH} - \mathbf{M}_{CE}|} \quad (25)$$

また、 \mathbf{V}_1 は右上腕、 \mathbf{V}_2 は右下腕のベクトルとし、同様にして求めた、 $\mathbf{V}_3, \mathbf{V}_4$ をそれぞれ左上腕、左下腕のベクトルとする。姿勢 \mathbf{P} は、単位方向ベクトル $\mathbf{V}_1, \mathbf{V}_2, \mathbf{V}_3, \mathbf{V}_4$ を用いて以下のように表す。

$$\mathbf{P} = [\mathbf{V}_1^T \ \mathbf{V}_2^T \ \mathbf{V}_3^T \ \mathbf{V}_4^T]^T \quad (26)$$

登録ジェスチャは式 (26) を用いて以下のように表す。

$$\mathbf{S}_g = [\mathbf{P}_{g1}^T \ \mathbf{P}_{g2}^T \ \dots \ \mathbf{P}_{gF}^T]^T \quad (27)$$

ここで、 \mathbf{P}_{gf} は第 g ジェスチャ \mathbf{S}_g 中の第 f フレームの姿勢を表す。

3.3.1 姿勢類似度

胴体の回転に不変な評価値を得るため、本稿では複素相関を用いて姿勢類似度を表す。世界座標における重力方向をカメラ座標の y 軸と一致させておくことで、 y 軸まわりの回転 (ほぼ胴体の回転に相当) に不変な類似度を計算できる。なお、 i を虚数単位とし、 \mathbf{P} の y 軸成分以外の 2 軸成分を用いた複素表現 \mathbf{q} (4 次元複素ベクトル) を次式で表す。

$$\mathbf{q} = \mathbf{P}_x + i \mathbf{P}_z \quad (28)$$

また、 \mathbf{q} の複素共役を \mathbf{q}^* で表す。

$$\mathbf{q}^* = \mathbf{P}_x - i \mathbf{P}_z \quad (29)$$

ここで、 $\mathbf{P}_x, \mathbf{P}_y, \mathbf{P}_z$ はそれぞれ \mathbf{P} の x, y, z 軸成分のみで構成される 4 次元ベクトルである。

2 つの姿勢 \mathbf{P}, \mathbf{P}' に対応する複素表現が \mathbf{q}, \mathbf{q}' で表せる場合、 \mathbf{q} と \mathbf{q}' の内積を以下のように表す。

$$\hat{\epsilon}_P = \mathbf{q}^* \mathbf{q}' \quad (30)$$

$$= |\hat{\epsilon}_P| e^{i\theta} \quad (31)$$

式 (31) で得られる $|\hat{\epsilon}_P|$ は、 y 軸まわりの回転に不変である。また、 θ は y 軸まわりの回転を示す。ここに、 Y 軸方向に関する項を加えることで、2 つの姿勢 \mathbf{P}, \mathbf{P}' の姿勢類似度 C を定義できる。

$$C(\mathbf{P}, \mathbf{P}') = \frac{1}{4} (|\hat{\epsilon}_P| + \mathbf{P}_y^T \mathbf{P}'_y) \quad (32)$$

この値は、 \mathbf{P}' を y 軸まわりに $-\theta$ だけ回転させた時の 2 姿勢間の正規化相関を示している。

3.3.2 ジェスチャ類似度

i を虚数単位とし、 \mathbf{S} の y 軸成分以外の 2 軸成分を用いた複素表現 \mathbf{Q} (120 次元複素ベクトル) を以下のように表す。

$$\mathbf{Q} = \mathbf{S}_x + i \mathbf{S}_z \quad (33)$$

また、 \mathbf{Q} の複素共役を \mathbf{Q}^* で表す。

$$\mathbf{Q}^* = \mathbf{S}_x - i \mathbf{S}_z \quad (34)$$

ここで、 $\mathbf{S}_x, \mathbf{S}_y, \mathbf{S}_z$ はそれぞれ \mathbf{S} の x, y, z 軸成分のみで構成される 120 次元ベクトルである。2 つのジェスチャ \mathbf{S}, \mathbf{S}' に対応する複素表現が \mathbf{Q}, \mathbf{Q}' で表せる場合、 \mathbf{Q} と \mathbf{Q}' の内積を以下のように表す。

$$\hat{\epsilon}_S = \mathbf{Q}^* \mathbf{Q}' \quad (35)$$

$$= |\hat{\epsilon}_S| e^{i\theta} \quad (36)$$

式 (36) で得られる $|\hat{\epsilon}_S|$ は、 y 軸まわりの回転に不変である。また、 θ は y 軸まわりの回転を示す。ここに、 y 軸方向に関する項を加えることで、2 つのジェスチャ \mathbf{S}, \mathbf{S}' の類似度 C' を定義できる。

$$C'(\mathbf{S}, \mathbf{S}') = \frac{1}{4F} (|\hat{\epsilon}_S| + \mathbf{S}_y^T \mathbf{S}'_y) \quad (37)$$

この値は、 \mathbf{S}' を y 軸まわりに $-\theta$ だけ回転させた時の 2 ジェスチャ間の正規化相関を示している。

3.4 ジェスチャの認識

3.4.1 ジェスチャ開始、終了検出

本稿のジェスチャ認識では、現在推定されている姿勢と各登録ジェスチャの初期フレームの姿勢と現在推定されている姿勢との類似度を求め、ジェスチャの開始を検出する。ここで、類似度は次式によって計算する。

$$\epsilon_g = C(\mathbf{P}_f, \mathbf{P}_{g1}) \quad (38)$$

これにより、胴体の回転に不変な評価値を得ることができ [8] [9]。また、 \mathbf{P}_f は、現在の姿勢推定結果、 \mathbf{P}_{g1} は \mathbf{S}_g の第 1 フレームの姿勢である。

各ジェスチャの開始が確認されると、ジェスチャの終了検出を行う。このとき用いられる類似度は次式で示される。

$$\epsilon'_g = C(\mathbf{P}_f, \mathbf{P}_{gF}) \quad (39)$$

ここで、 \mathbf{P}_{gF} は第 g ジェスチャ \mathbf{S}_g の第 $F (= 30)$ フレームの姿勢である。 $\epsilon_g > \tau_1, \epsilon'_g > \tau_1$ のときに開始、終了を検出する。ここで、 τ_1 は閾値である。

3.4.2 ジェスチャの対応付け

ジェスチャを行う場合、そのときどきで、全体、または部分的な速さが登録されているジェスチャと異なることがある。そこで、Algorithm 1 で示す動的計画法を用いて、最適な入力ジェスチャ $\hat{\mathbf{S}}$ を求める (ジェスチャの開始フレームから終了フレームまでの推定姿勢が登録姿勢のどのフレームと対応するか求める)。ここで、 $C_{i,j}$ は、登録ジェスチャの第 i フレームと入力画像の第 j フレームの類似度を示す。

3.4.3 評価

得られた $\hat{\mathbf{S}}$ を用いて式 (40) により評価を行う。

$$\epsilon_{S_g} = C'(\hat{\mathbf{S}}, \mathbf{S}_g) \quad (40)$$

$\epsilon_{S_g} > \tau_2$ のときジェスチャ g を認識したとする。ここで、 τ_2 は閾値である。

Algorithm 1 for \tilde{S}

```

 $i \leftarrow I, j \leftarrow J$ 
while  $i \neq 1 \vee j \neq 1$  do
  if  $\max\{C_{i-1,j}, C_{i,j-1}, C_{i-1,j-1}\} = C_{i-1,j}$  then
     $(i, j) \leftarrow (i-1, j)$ 
  else if  $\max\{C_{i-1,j}, C_{i,j-1}, C_{i-1,j-1}\} = C_{i,j-1}$  then
     $(i, j) \leftarrow (i, j-1)$ 
  else
     $(i, j) \leftarrow (i-1, j-1)$ 
  end if
end while

```

3.5 多段階認識による最適化

3.4節で述べた手法について、登録ジェスチャを30フレームとして扱う場合だけでなく、 L 段階に分け、 $30/L$ フレームとして扱った場合を考える。ここで、第 l ($l = 2, 3, \dots, L$)段階における開始検出は、第 $l-1$ 段階の終了検出とする。第 L 段階目の式(40)における ϵ'_g が閾値を越えた場合そのジェスチャを認識したとする。これに伴い、第1段階目の終了検出は式(39)でなく、次の式による。

$$\epsilon_{gl} = C(P_f, P_{g(FI/L)}) \quad (41)$$

4. 実験

4.1 ジェスチャの種類

述べた追跡、認識手法について7種類のジェスチャを登録して実験を行った。各登録ジェスチャの初期フレームから最終フレームの一部を図4に図示する。各ジェスチャにつき30個のシーケンスを用意した。ジェスチャ1は右手を右上から左上へ動かすジェスチャである。ジェスチャ2は胴体の前で両手をあわせるジェスチャであり、胴体と下腕でオクルージョンが発生する。ジェスチャ3は右腕を右から左に動かすジェスチャであり、比較的短いジェスチャである。ジェスチャ4は胴体の横で両腕を下から上へ動かす、さらに下へ動かすジェスチャであり、比較的長いジェスチャである。ジェスチャ5, 6は右腕を胴体の前で回すジェスチャであり、それぞれ回転する方向が逆になっている。ジェスチャ7は胴体の前で両腕を回すジェスチャであり、オクルージョンが複雑で、難しく設定したジェスチャである。学習では、各ジェスチャについて15シーケンスの追跡結果から平均ジェスチャを作成した。なお、学習においては、ジェスチャの開始、終了は手動で与えている。学習に用いなかった残り15シーケンスをテストシーケンスとし、ジェスチャ認識の実験を行った。

姿勢推定により、各ジェスチャの初期姿勢を検出したときに、それぞれのジェスチャで認識を開始する。開始フレームがほぼ同じ姿勢であるジェスチャが存在する場合、ジェスチャ認識は並列に実行されることになる。

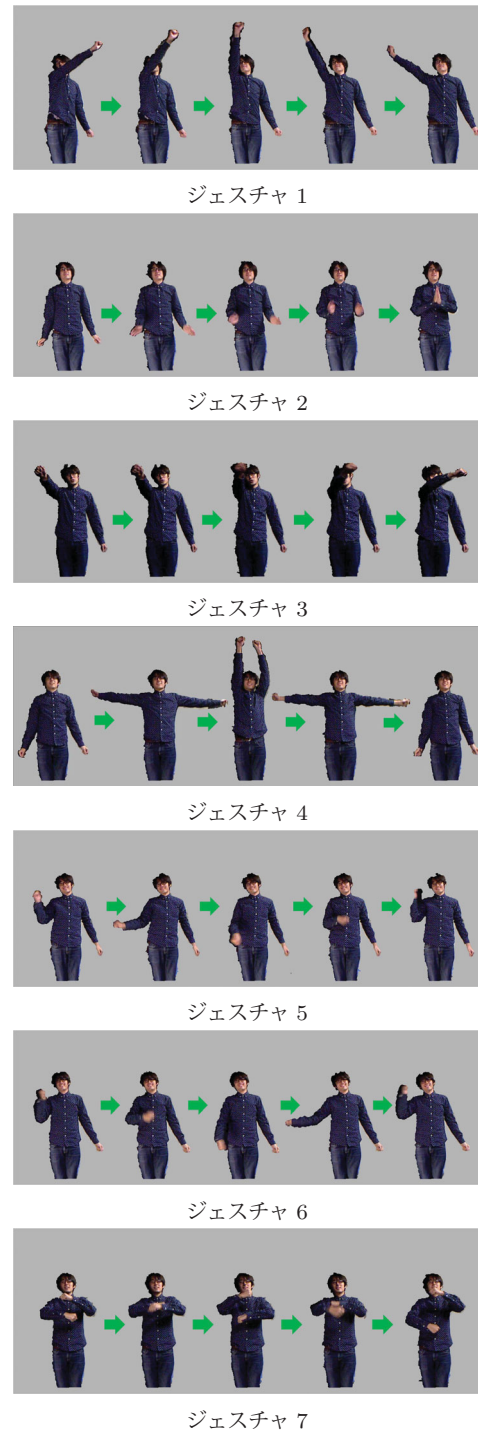


図4 登録ジェスチャ

4.2 実験条件

本稿では、距離・画像センサとして、Kinect for Xbox360を用いた。このセンサは赤外線パターンを前方に照射し、赤外線カメラで取得した画像から距離を計算することで、距離情報をリアルタイムで獲得できる。距離情報はmm単位で表され、50cm以上離れた対象について有効な情報が得られる。

追跡に用いた人物モデルを図5に示す。これらのモデルは胴体と左右上腕、左右下腕からなり、それぞれKinectで取得したデータから作成した。

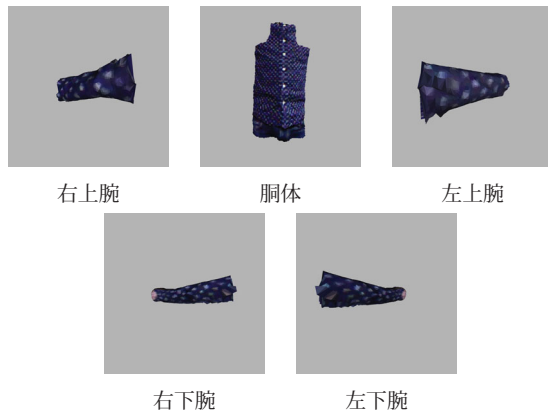


図 5 上半身モデル

表 1 追跡条件

項目	設定値
サンプル数	1500
粗密探索の段数	3
初段のサンプル生成範囲 (σ)	表 2 参照
疎点の数	128
見え画像と奥行き情報の比率 (λ)	0.90

表 2 粗密追跡の初段のサンプル生成に用いた標準偏差
($t_1 \sim t_3$ は mm, その他は deg 単位である)

胴体						肩		肘	
θ_0	ϕ_0	ψ_0	t_1	t_2	t_3	ϕ_S	θ_S	ϕ_E	θ_E
1.0	1.0	1.0	10.0	0.5	7.0	8.0	8.0	12.0	12.0

追跡系は 3 段階からなる粗密追跡法 [7] をベースとして作成し, 14 次元姿勢空間を取り扱っているパラメータ設定の詳細を表 1, 2 に示す. なお, ジェスチャの開始, 終了に用いる閾値 τ_1 , ジェスチャ認識に用いる閾値 τ_2 はともに 0.96 と設定した.

4.3 複素相関の有効性の検証

本稿で示した手法は, カメラ座標のベクトルで表現したジェスチャを複素相関により評価を行うものである. これに対し, 胴体座標で表現したベクトルを実相関により評価する比較実験も行い, その認識精度を確認する. なお, 実相関では, 式 (40) のかわりに式 (42) で得られる相関が閾値を越えた時ジェスチャ g を認識する.

$$C'_R(\mathbf{S}, \mathbf{S}') = \frac{\mathbf{S}^T \mathbf{S}'}{\|\mathbf{S}\| \|\mathbf{S}'\|} \quad (42)$$

胴体座標で表現することにより, 胴体の姿勢推定が正確に行える場合には, 実相関によってもジェスチャ認識を行える. 複素相関を用いた場合との比較を表 3 に示す. 実相関を用いた場合, 認識率が低くなっているのは, 胴体, 上腕の姿勢推定誤差の影響と考えられる. これに対し, 複素相関を用いる場合には, 姿勢推定誤差の影響をあまり受けず, いずれのジェスチャにおいても高い認識率が得られている.

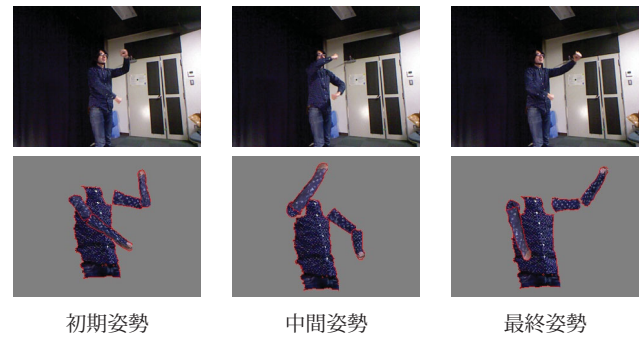


図 6 向きの違う入力画像

表 3 実相関と複素相関 (多段階) による認識率 (%)

	実相関	複素相関			
	1 段階	1 段階	2 段階	5 段階	10 段階
ジェスチャ 1	80%	94%	94%	94%	94%
ジェスチャ 2	100%	100%	100%	100%	60%
ジェスチャ 3	87%	100%	100%	100%	100%
ジェスチャ 4	94%	100%	100%	94%	100%
ジェスチャ 5	100%	100%	100%	100%	94%
ジェスチャ 6	100%	100%	100%	100%	100%
ジェスチャ 7	80%	100%	100%	100%	100%

次に, 正面から約 15 度, 30 度, 45 度の各方向を向いて同じジェスチャを行った場合についても実験を行った. この結果, 実相関を用いた場合には, 約 15 度までの簡単なジェスチャしか認識できなかったのに対し, 複素相関を用いた場合には, 45 度までのいずれのジェスチャにおいて, 正面向き同様に精度の良い識別ができた. 図 6 6 に入力動画の一部と姿勢推定結果を示す. なお, この例は正面から約 45 度の方向を向いており, ジェスチャ 7 を行っているシーケンスである. 以上の結果から, 本稿で示した複素相関によるジェスチャ認識の有効性を確認することができた.

4.4 多段階認識による比較

複素相関を用いる方法について, 3.5 節で述べた多段階認識に関する比較実験を行った. なお, 多段階認識では, 第 L 段階目の認識が確認できた場合にジェスチャ認識が成功したものとす. 実験結果を表 3 に示す.

1, 2 段階認識では, すべてのジェスチャにおいて高い認識率となった. それに対し, 10 段階認識の場合, ジェスチャ 2 の認識率が 60% となった. ジェスチャ 2 は短いジェスチャであり, 常にオクルージョンが発生しているため推定姿勢に多くのノイズが加わっている. 10 段階認識の場合, 3 フレーム単位で段階的に認識していくため, 短時間の姿勢推定誤差によって段階の進行に乱れが生じ, 認識率が低くなったと考えられる.

また, 提案手法では, ジェスチャの開始, 終了が検出された時点でジェスチャ類似度の評価を行う. そのため, ジェスチャの数が多く, 段階数が少ない場合, たまたま似た姿勢が数フレーム入力されると誤認識を行う可能性がある. こ

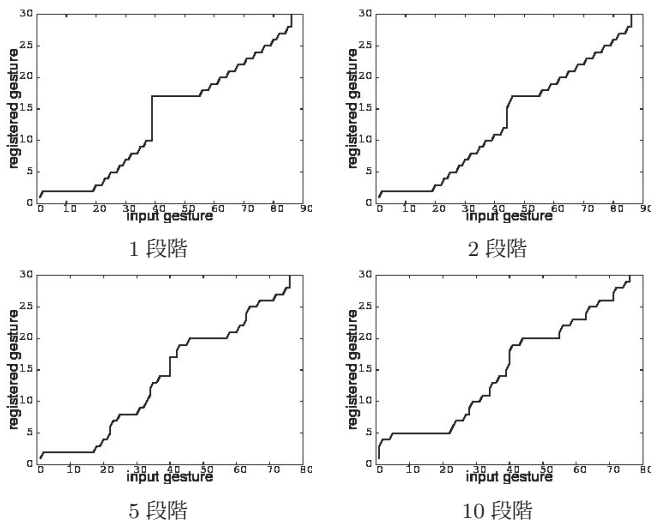


図 7 入力ジェスチャと登録ジェスチャ 4 の対応

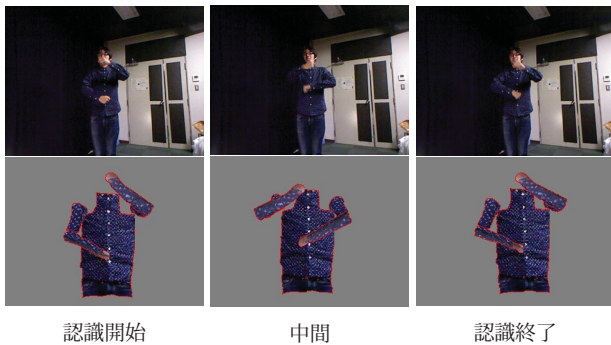


図 8 入力と推定結果 ジェスチャ 7

のため、適切に部分ジェスチャへの分割を行うことが望ましい。

図 7 に、ジェスチャ 4 について、異なる段階数における登録ジェスチャと入力ジェスチャとの対応付け結果を示す。この図で、縦に一直線になっている場所は入力における 1 フレームが登録の複数フレームに対応していることになる。逆に、横に一直線になっている場所は、複数の入力が登録の 1 フレームに対応している。この場合、ジェスチャ認識には複数ある入力の内もっとも評価値が大きい入力フレームのみが用いられることになる。ジェスチャ 4 には、ジェスチャの途中で動きが止まる場所が存在するため、登録ジェスチャと入力ジェスチャがうまく対応していないように見えるが、この間は姿勢が変わっていないためジェスチャ認識としては高い評価値を得ることができ、正しくジェスチャを認識できている。また、実験で使用したジェスチャ 7 の認識開始フレームから終了フレームまでの一部を、それらに対する姿勢推定結果とともに図 8 に示す。

4.5 考察

本稿で提案した手法では、人物の胴体座標基準にしたジェスチャ認識を取り扱ったが、人間のジェスチャには、本手法では特定できないものもある。例えば、「指をさす」と

いう動作は、世界座標を用いないと解析できない。また、より直感的なジェスチャ認識を目指す場合、左右の腕でジェスチャを独立させる必要があると考えられる。また、今回示した追跡・認識系はリアルタイム処理で実現できていない。認識ステップでは、ジェスチャの認識が開始されると 1 つにつき約 0.1 msec である。一方、現在の 1 コア CPU 上での実装では、追跡ステップにおいて 1 フレームにつき約 400 msec を要する。今後、プログラムの最適化や GPU 実装により実時間処理の実現を目指す。

5. まとめ

本稿では距離・画像センサを用いた人物追跡によるジェスチャ認識について議論した。追跡時に得られる姿勢空間パラメータを時系列として保持し、ジェスチャとして登録し、重力方向まわりの回転に不変な類似度を用いて評価を行った。これにより、人物の向きが変わっても、正面を向いている時同様、ジェスチャの認識が行えることを示した。また、登録ジェスチャを分割し、それぞれで入力を段階的に評価していくことによる認識への影響についても述べた。今後の課題として、ジェスチャの登録数の増加、初期姿勢の検出、複雑なオクルージョンへの処理、高速化が挙げられる。

参考文献

- [1] A.F.Bobick and J.W.Davis, "Real-time recognition of activity using temporal templates," Proc. 3rd IEEE Workshop on Applications of Computer Vision (WACV), pp.39-42, 1996.
- [2] Y. Song, D. Demirdjian and R. Davis, "Tracking body and hands for gesture recognition: Natops aircraft handling signals database," Proc. 9th IEEE Conference on Automatic Face and Gesture Recognition (FG 2011), pp.500-506, March 2011.
- [3] T. Shakunaga, "Pose estimation of jointed structures," Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR'91), pp.566-572, 1991.
- [4] 佐竹純二, 尺長健, "階層的注視点制御による動画像上での複数人物追跡," 電子情報通信学会論文誌 (D-II), no.8, pp.1212-1221, 2003.
- [5] 四宮洋平, 加藤秀章, 尺長健, "Kinect を用いた 3 次元疎テンプレートによる人物姿勢追跡," 電子情報通信学会技術研究報告 PRMU2011-264, 2012.
- [6] 松原康晴, 尺長健, "疎テンプレートマッチングとその実時間物体追跡への応用," 情報処理学会論文誌, vol.46, no.SIG(CVIM11), pp.60-71, 2005.
- [7] Y. Oka, T. Kuroda, T. Migita and T. Shakunaga, "Tracking 3d pose of rigid object by sparse template matching," Proc. International Conference on Image and Graphics, pp.390-397, 2009.
- [8] 山根亮, 川嶋幸治, 戸高千智, 尺長健, "動作データの時系列相関行列による舞踊動作解析," 電子情報通信学会論文誌 (D-II), no.8, pp.1652-1661, 2005.
- [9] R. Yamane and T. Shakunaga, "Dance motion analysis by correlation matrix between pose sequences," Proc. International Conference on Image and Vision Computing, New Zealand (IVCNZ2010), 2010.