人検出結果に基づくカメラの自己キャリブレーションと3次元位置推定

安藤寛哲,藤吉弘亘(中部大学)

A Method of Measuring 3D Position and Camera Self-Calibration Based on Results of Human Detection Hiroaki Ando, Hironobu Fujiyoshi (Chubu University)

Abstract

A camera self-calibration method based on the results of human detection is proposed. This method extracts the positions and heights of people in the target scene from the results of human detection and human-area segmentation, and estimates camera parameters such as a location of camera in the world coordinate and a vanishing line in the image coordinate. Calibrating a camera generally requires intensive work, but the proposed technique can perform self-calibration using parameters that are automatically extracted from the target image. As a result, our method can estimate the three-dimensional position of an object even with a camera that had not been previously calibrated. Experimental results show that the accuracy of estimated camera parameters can be improved by using the results of human-area segmentation.

キーワード:自己カメラキャリブレーション,人検出,セグメンテーション,3次元位置推定 (Camera self-calibration, Pedestrian detection, Segmentation, 3D position estimation)

1. はじめに

近年,公共施設だけではなくオフィスや一般家庭への監 視カメラの普及とともに,物体検出や追跡などの動画像処 理技術による監視システムの需要が高まっている.特に,映 像中から人を検出⁽¹⁾⁽²⁾し、人の3次元位置情報や身長を 知ることは,場所に応じた人数カウントや人流測定におい て重要である.小川らは,人の身長と位置情報の推定に,3 次元空間の光線情報に基づく手法を提案した⁽³⁾.この手法 では, Tsai モデル⁽⁴⁾ に基づきカメラキャリブレーションを 行い,カメラの内部パラメータと外部パラメータを用いて 人の位置と身長の推定を行う.キャリブレーションにより 求めたカメラパラメータを利用することで,画像中の2次 元座標と世界座標における3次元座標の関係を求めること が可能となる.しかし,一般的にカメラキャリブレーショ ンは,画像中の座標に対応する世界座標を与える必要があ り,広い範囲でのキャリブレーションは難しく,正確に行 うには大変な労力が必要となる.

一方,手間を必要としない自己カメラキャリブレーショ ンの手法も多く提案されている.Lvらは,歩行者を背景モ デルとの差分から算出し,歩行者の頭と足下の位置検出し, 画像中の人の高さ世界座標と画像座標の関係を用いて,消 失点と水平線を算出することにより,カメラの自己キャリ プレーション手法を提案している⁽⁵⁾.しかし,Lvらの手 法では,ノイズ等の影響により画像中から人の大きさを安 定して取得できないという問題がある.また,Hoiemらは, 入力画像から物体の検出を行い,検出結果から得られた複 数の物体高さの情報と,入力画像における三次元の面構造 と消失点を推定して,自己カメラキャリブレーションする 手法を提案している⁽⁶⁾.これらの物体検出に基づくカメラ の自己キャリブレーション手法では,参照する物体の高さ (例えば人の高さ)を精度良く求める必要がある.

そこで,本稿では人領域のセグメンテーション結果に基 づくカメラの自己カメラキャリブレーション手法を提案す る.提案手法は,推定するシーンに対し人検出と人領域の セグメンテーションを行う画像から抽出した複数の人の高 さと位置情報と透視投影モデルにおける世界座標と画像座 標の関係を用いてカメラの自己キャリブレーションを行う. また,自己キャリブレーションの結果を用いることにより, 未校正のカメラにおいても3次元位置の情報を推定するこ とができる.

2. 人領域のセグメンテーション

提案手法では,単眼のカメラで撮影された映像から,画像中の物体の高さと位置,透視投影モデルにおける世界座標と画像座標の関係を用いて,カメラの自己キャリブレーションを実現する.本手法では,画像中の物体の高さ情報を用いてカメラパラメータを高精度に推定するため,人領域を正確にセグメンテーションする必要がある.

本手法では,人領域のセグメンテーション手法に,村井 らの提案した弱識別器の応答による類似シルエットの選択 を用いたセグメンテーション手法⁽⁷⁾を用いる.この手法は, 人検出器の学習と学習サンプルに対応するシルエット画像 のスコアからハッシュテーブルを求めるオフライン処理と, 入力画像からラスタスキャンによる人検出を行い Chamfer Matchingを用いてセグメンテーションを行うオンライン処 理から構成される.処理の流れを図1に示し,以下に詳細 を述べる.

2・1 人検出器の構築 人検出法として HOG (Histograms of Oriented Gradients)特徴量⁽¹⁾を用いた Real AdaBoost⁽⁸⁾による人検出器を用いる.HOG 特徴量は1つの局所領域内におけるエッジ方向ごとのエッジ強度に着目



図 1 人領域のセグメンテーション処理の流れ Fig. 1. Processing flow in the human-area segmentation

した特徴量であり,照明変動と,局所的な幾何学的変化に 頑健な特徴量である.人検出では,このHOG特徴量を予 め用意した学習サンプルから算出し,Real AdaBoost によ り人検出器の構築を行う.

2・2 シルエット画像のスコア算出 入力画像から 人検出を行い,シルエット画像のマッチングを行う候補領 域を推定する.構築した検出器を用いてシルエット画像の スコアを算出する.シルエットのスコアとは,人検出器の 各弱識別器の応答値を多次元のベクトルとして表現したも のである.まず,学習サンプルからシルエット画像を作成 する.次に,作成したシルエット画像に対応する学習サン プル(ポジティブのみ)を構築した人検出器に入力する.入 力画像に対して各弱識別器は人か人以外を識別した結果と して返す.T個の弱識別器に対応したT次元の特徴ベクト ルを対応するシルエット画像のスコアとする.

2・3 ハッシュテーブルの作成 算出したシルエット 画像のスコアは、Chamfer Matching を行う対象領域に類似 したシルエットを選択するために使用する.シルエット画像 のスコアを記録したハッシュテーブルを作成することで高 速な類似シルエット画像検索を実現する.ある一つのシル エット画像 y から算出された弱識別器の応答に基づくスコ アを、T 次元の特徴ベクトル $y = \{h_1(y), h_2(y), \dots, h_T(y)\}$ とする.まず、次式により各次元の特徴量を 2 値化し、ビッ トベクトル $u = (u_1, u_2, \dots, u_T)$ を作成する.

次に,以下の式を用いて,ハッシュテーブルのインデックス *H_{index}* を求め,ハッシュテーブルに特徴ベクトルとシル エット画像の ID を登録する.

$$H_{index} = \left(\sum_{i=1}^{T} u_i 2^i\right) \mod H_{size} \cdots \cdots \cdots \cdots (2)$$

ここで, *H_{size}* はハッシュテーブルのサイズである.以上の処理をすべてのシルエット画像に対して行い, ハッシュ テーブルを作成する.

2.4 類似シルエットの選択 オンライン処理として 実際の入力画像から人検出を行い, MeanShift クラスタリ ングによりウインドウを統合する.検出されたウインドウ 内の領域が検出器に入力されたときの各弱識別器の応答値 のT次元の特徴ベクトルとシルエット画像のスコアとして 記録した特徴ベクトルを比較することで類似シルエット画 像の選択を行う.その方法として,作成したハッシュテー ブルを用いて近似最近傍探索を行う.人検出により検出さ れた領域から得られた弱識別器の応答に基づく特徴ベクト ルを $\mathbf{x} = (h_1(x), h_2(x), \cdots, h_T(x))$ とする.この特徴ベク トルに対しても式(1),(2)を用いてハッシュのインデック スを求め,同じハッシュのインデックスを持つ特徴ベクト ルをハッシュテーブルから探索することで類似する特徴べ クトルを持つシルエット画像を選択する.

2.5 Chamfer Matching によるセグメンテーション 人検出により検出された領域と,弱識別器の応答値から選 択された類似シルエット画像を Chamfer Matching⁽⁹⁾ によ リマッチングすることで人領域のセグメンテーションを行 う.入力画像として人検出により検出された領域,テンプ レート画像としてシルエット画像を用いて相違度に基づい てマッチングを行う.このとき,人検出により検出された 領域は必ずしもきれいに人領域を検出しているとは限らな いため, Chamfer Matching を行う領域としてマージンを 持たせて少し大きめに切り出した領域を使用する.また, Chamfer Matching はスケール変化に敏感であるため,シ ルエット画像のスケールを変化させながらマッチングを行 う. 最終的に, 最も Chamfer Matching の相違度が小さく なる位置とスケールとシルエット画像を決定し,シルエッ ト画像が重なっている部分を切り出すことで人領域のセグ

メンテーションを行う. セグメンテーション結果から,画像中の人の高さ h_i と位置 v_i を精度よく抽出することができる.

3. カメラの自己キャリブレーション

人領域のセグメンテーション結果から得られた画像中の 人の高さを用いてカメラの自己キャリブレーションを行う. まず,文献⁽⁶⁾を参考に,透視投影モデルにおけるカメラ位 置と人の高さと人の位置の関係を定式化する.本章では,画 像座標と世界座標とカメラパラメータの関係と,その関係 からカメラパラメータの導出方法について述べる.

3·1 カメラ位置と人の高さの関係 図2に,カメ ラ位置と人の高さの関係を示す.



図 2 カメラと人の高さの関係 Fig. 2. Relationship between camera and human height

画像の左下の座標を (0,0) とし,画像の縦と横を画像の 縦幅のサイズで正規化した画像の座標系を (u,v) で与える. また,世界座標系を (x,y,z) とする.ただし, y を高さ, z を奥行きとする.次に,カメラのチルト角を θ ,焦点距離 を f,カメラ中心 (u_c, v_c) ,カメラの高さを y_c とする.世界 座標系は,カメラ位置を基準として, $z_c = 0, x_c = 0$ とし, 接地平面を y = 0 と定義する.また,カメラ中心 (u_c, v_c) は画像中心とし,焦点距離 f は 1.4 とする.カメラにロー ルは発生せず,画像座標系の地面の消失ラインとして水平 線を v_0 と定義する.このとき,カメラのチルト角(ラジア ン) は次式で与えられる.

カメラモデルとして, 歪み無しの単位アスペクト比である 透視投影モデル(Perspective projection model)を使用す ると,世界座標系から画像座標系への変換は以下の式で与 えられる.

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \frac{1}{z} \begin{bmatrix} f & 0 & u_c \\ 0 & f & v_c \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos\theta & -\sin\theta & y_c \\ 0 & \sin\theta & \cos\theta & 0 \end{bmatrix} \begin{vmatrix} x \\ y \\ z \\ 1 \end{vmatrix}$$
(4)

3・2 カメラパラメータの導出 世界座標系と画像 座標系とカメラパラメータの関係から,カメラパラメータ の導出方法について説明する.式(4)を物体の高さ y につ いて解くと次式となる. ここで,画像中の人(直立と仮定)領域の上底と下底の位置をそれぞれ $v_t \ge v_b \ge 0$ て与える(図2).下底位置 v_b で地面と接する.つまり v_b のときy = 0なので,物体の奥行きzは次式のように求められる.

 $z = fy_c / (f \sin \theta - (v_c - v_b) \cos \theta) \cdots \cdots \cdots \cdots \cdots (6)$

式(5)と式(6)から,物体の高さyは次式のようになる.

$$y = \frac{fy_c(f\sin\theta - (v_c - v_t)\cos\theta)/(f\sin\theta - (v_c - v_b)\cos\theta) - fy_c}{(v_c - v_t)\sin\theta + f\cos\theta}$$

(7)

ここで,カメラのチルト角が小さい(推定された水平線の位置が画像内に存在する)場合, $\cos \theta \approx 1$, $\sin \theta \approx \theta$, $\theta \approx \frac{v_c - v_0}{f}$ と近似できる.従って,式(7)は以下のように 表現できる.

$$y \approx y_c \frac{v_t - v_b}{v_0 - v_b} / \left(1 + (v_c - v_0)(v_c - v_t) / f^2 \right) \dots \dots (8)$$

分母の $(v_c - v_0)(v_c - v_t)/f^2$ は, チルト角が小さいとき $v_c - v_0 \approx 0$ となるため, $(v_c - v_0)(v_c - v_t)/f^2 \approx 0$ と近似 できる. 従って,式 (8) は以下のようになる.

ここで,画像中に存在する i 番目の人領域を考えたとき, $v_t - v_b$ は画像中の i 番目の人領域の高さ h_i とする.また, 人領域の基準座標を矩形領域の左下に設定したとき, $v_b = v_i$ となる.従って,求めたい式は以下のようになる.

式 (10) を変形し, n 人のサンプルとの関係を行列で表現す ると次式のようになる.

式 (11) からカメラの高さ y_c と水平線の位置 v₀ を求める.

3・3 RANSAC によるカメラパラメータ推定 式 (11)において,既知のパラメータは,2・5 で述べたように, 人領域のセグメンテーション結果より画像中から抽出した 人の高さ h_i と人領域の位置 v_i である.世界座標における カメラの高さ y_c と画像座標上の水平線の位置 v₀ を求める ためには,世界座標系での人の高さ y_i が必要である.しか し,与えられている映像からそれぞれの世界座標系の人の 高さの真値を知ることは不可能である.そこで,提案手法 では人の身長として平均値 170.0[cm],標準偏差 8.0 の正規 乱数により人の高さを与え,RANSAC(RANdam SAmple Consensus)⁽¹⁰⁾を用いてカメラの高さとチルト角の推定を 行う.RANSAC の流れは以下の通りである. Step1 ランダムに2点のサンプルを選択

Step2 最小二乗法によりカメラの高さと水平線位置を推定 Step3 推定値から全サンプルとの誤差を算出

Step4 誤差の中間値を選択

Step5 Step1 ~ Step4 を繰り返し,誤差の中間値が最 小となるパラメータを求め,最終的なカメラの高さと水 平線位置とする

3・4 人の3次元位置推定 推定したカメラパラメー タと画像中における物体の位置 (u,v)から,世界座標にお ける物体の3次元位置 (x,y,z)を推定する.物体の奥行き zは,世界座標系と画像座標系の変換式から,y = 0のとき 求めたカメラパラメータのカメラの高さと,画像座標にお ける位置関係から次式となる.

$$z = fy_c / (f \sin \theta - (v_c - v_b) \cos \theta) \cdots \cdots \cdots \cdots \cdots (12)$$

次に,式(12)から求めた奥行き *z* と画像中の位置から,*x* と *y* を求めることが出来る.世界座標における高さ *y* は, 奥行き *z* と,カメラの高さ *yc* と画像中での高さの関係から 次式となる.

世界座標における *x* の推定には,世界座標系から画像座 標系に変換する式 (4) を画像座標 *u* について解くと次式と なる.

式 (14) から, x は次式のようになる.

式 (15) から求めた世界座標における x 座標と,式 (13) と 式 (12) から求まる y, z 座標の値を,自己カメラキャリプ レーションで求めたカメラパラメータを用いて算出するこ とで,対象物体の 3 次元位置を推定することができる.

4. 実 験

4・1 カメラパラメータの推定

4・1・1 実験概要 本実験では,撮影した映像から人 領域のセグメンテーションを行い,その結果から撮影した カメラの自己キャリブレーションを行う.求めるパラメー タはカメラの高さ y_c と水平線位置 v₀ で,水平線位置 v₀ と 式(3)からカメラの角度 θ である.また,推定したパラメー タを用いて現実世界での3次元位置の推定を行う.撮影し た映像は,640×360 ピクセル,推定に用いたセグメンテー ション結果のサンプル数は 600 である.

4・1・2 実験結果 図 3(a), (b) にカメラパラメータ の推定に用いた 600 サンプルの人検出結果の分布を示す.表 1 にパラメータ推定結果を示す.表1より,セグメンテー ションありの場合にカメラ高さが真値との誤差が約0.3 [cm], カメラの角度が真値との誤差は約3.3 [degree] の精度で推



(c) Example of estimation

図 3 パラメータ推定結果 1(オリジナルシーケンス)
$$y_c = 190.93$$
[cm], $v_0 = 0.85, \theta = 14.4$ [deg]
Fig. 3. Estimation result 1(Original Sequence)

定できた.図3(c),図4に自己キャリブレーションにより 推定した水平線位置 v_0 , y_c と式(10)から推定した検出し た人の身長, θ と式(12)から求めた奥行きzの値を示す. 提案するカメラの自己キャリブレーション法は,カメラの 内部パラメータを推定することはできないが,推定したカ メラの高さ y_c と水平線位置 v_0 を用いることで人検出の誤 検出抑制やトラッキング精度の向上が期待できる.

表 1 カメラパラメータ推定結果

Table 1. Estimation results of camera parameter

Camera parameter	True value	Segmentation	
		Proposed	Conventional
Camera height yc[cm]	184.6	184.3	214.1
Camera angle θ [degree]	10.0	13.3	18.1

4・2 人の3次元位置の推定 推定したカメラパラ メータを用いて,検出した人の世界座標における3次元位 置を推定をする.推定精度を世界座標の真値と比較するこ とで,推定精度の評価を行う.図5に示す15点のランド マーク点の画像座標を用いて3次元位置の精度評価を行う.

4・2・1 評価実験 図6に3次元位置の推定結果を 示す.また,表2に3次元位置推定誤差のx軸,z軸の平 均と分散を示す.セグメンテーションした結果により推定 したカメラパラメータを用いて3次元位置を推定すること で,推定誤差を小さくすることができた.また,図6の結 果から提案手法で3次元位置をおおまかに推定できている ことから,数mの領域に人がどのように存在しているかな どの位置関係を推定できることが分かる.



図 4 パラメータ推定結果 2(CAVIAR dataset 2) $y_c = 350.05$ [cm], $v_0 = 1.12, \theta = 29.2$ [deg] Fig. 4. Estimation result 2(CAVIAR dataset 2)



図 5 推定する 3 次元位置 (オリジナルシーケンス) Fig. 5. Estimation for 3D positions (Original Sequence)



図 6 3 次元位置の推定誤差 (オリジナルシーケンス) Fig. 6. Estimation error of 3D positions(Original Sequence)

図7,図8に提案手法を用いた3次元位置の推定結果を 示す.未校正のカメラから自己キャリプレーションを行い, その結果を用いることで検出した人の3次元の位置関係や 身長を推定することが可能となる.

表 2 3 次元位置の誤差の平均と分散 [cm]

Table 2. Average and variance of error in 3D positions [cm]

	With Segmentation		Without Segmentation		
	Average	Variance	Average	Variance	
x-axis	5.6	14.1	9.1	76.1	
z-axis	20.6	43.9	33.8	573.0	

5. おわりに

本稿では、人検出とセグメンテーション結果に基づく、カ メラの自己キャリプレーション手法を提案した.セグメン テーションの結果から画像中の人の高さと、位置を自動的 に抽出することで自己キャリプレーションを行い、カメラ パラメータの推定を行った.これにより、推定シーンにお ける物体の3次元位置推定が可能となり、評価実験により その有効性を確認した.

今後は推定した物体の3次元位置情報を用いた高精度な 人の追跡を行う予定である.

文 献

- (1) N. Dalal and B. Triggs, "Histograms of Oriented Gradients for Human Detection", IEEE Computer Vision and Pattern Recognition, vol.1, pp.886–893, 2005.
- (2) Bo Wu and Ram Nevatia, "Detection and Tracking of Multiple, Partially Occluded Humans by Bayesian Combination of Edgelet based Part Detectors", International Journal of Computer Vision, 75(2), pp.247– 226, 2007.
- (3) 小川雄三,藤吉弘亘,"実空間に対応した Master-Slaving による追尾カメラシステム",第9回画像センシングシン ポジウム, June, 2003.
- (4) R. Y. Tsai: "A Versatile Camera Calibration Technique for High-Accuracy 3D Machine Vision Metrology Using Off-the-Shelf TV Cameras and Lenses", IEEE journal of Robotics and Automation, Vol. RA-3, No.4, pp. 323–344, 1999.
- (5) F. Lv, T. Zhao, R,Nevatia "Camera calibration from video of a walking human ", IEEE Trans. on Pattern Analysis and Machine Intelligence, vol.28, no.9, pp.1513-1518, September 2006.
- (6) D. Hoiem, A. A. Efros and M. Hebert, "Putting Objects in Perspective", International Journal of Computer Vision, vol. 80, No. 1, 2008.
- (7) 村井陽介,藤吉弘亘,"弱識別器の応答に基づく類似シル エット画像選択による Chamfer Matching を用いた人領 域のセグメンテーション",第12回画像の認識・理解シン ポジウム,2009.
- (8) R. E. Schapire and Y. Singer, "Improved Boosting Algorithms Using Confidence-rated Predictions", Machine Learning, No. 37, pp. 297–336, 1999.
- (9) D.M. Gavrila, "Multi-feature Hierarchical Template Matching Using Distance Transforms", IEEE International Conference on Pattern Recognition, pp. 439–444, 1998.
- (10) M. Fischer, and R. Bolles, "Randam sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography", communications of the ACM 24, pp. 381-385, 1981.







図 8 3 次元世直の推足結果 2(CAVIAR dataset 2) $y_c = 350.05$ [cm], $v_0 = 1.12, \theta = 29.2$ [deg] Fig. 8. Estimated results of 3D positions 2(CAVIAR dataset 2)