

アフィンカメラ補正による三次元位置の線形推定法

木下 敬介

ATR 人間情報通信研究所

張 正友

ATR 人間情報通信研究所

/ INRIA Sophia-Antipolis

複数のカメラで撮影した画像から、三次元空間中の点の三次元位置を、簡単に、そして、安定に推定する手法を提案する。従来の透視変換モデルを使うかわりに、アフィンカメラモデルを採用した。まず、複数のアフィンカメラで撮影された画像から、三次元空間の点の位置を線形な計算で推定できることを示す。次に、仮想的なアフィンカメラで撮影したのと等価になるように、現実のカメラで得られた画像を補正する方法を提案する。補正後の画像からは、三次元位置を線形な演算だけで推定することができる。

3D Reconstruction by Affine Camera Compensation

Keisuke Kinoshita

ATR Human Information
Processing Research Labs.

Zhengyou Zhang

ATR Human Information
Processing Research Labs.
/ INRIA Sophia-Antipolis

This paper proposes a simple and robust method that estimates 3D position of the points from multiple cameras. Affine camera model is used instead of conventional perspective camera model. First we show that, using multiple cameras, the 3D position is estimated from the image coordinates by linear operation. Then, we assume that Affine cameras are virtually placed in the 3D space. The image coordinates of the real cameras are modified to coincide with those of the virtual Affine cameras. 3D positions of the points are linearly estimated from this modified image coordinates.

1 はじめに

カメラから得られた画像をもとに、ロボットを制御したりする場合、環境の3次元情報が必要になる。従来から、運動視、ステレオ視などにより、3次元情報を復元する手法が様々開発されてきた。これらの手法は幾何学的な画像間、そして、画像と3次元空間の幾何学的な関係を利用したものである。これらは、透視変換という非線形な関係でモデル化している。そのため、画像から3次元を復元するという操作は非線形な導出となる。

それに反し、線形な関係だけで画像と3次元空間の関係を記述する方法もある。アフィンカメラと呼ばれるカメラモデルがそれである。これは、weak-perspective, para-perspectiveなどのカメラモデルを包括するモデルである。アフィンカメラモデルが仮定できるならば、3次元から画像、また、画像から3次元へはすべて線形な形で記述できる。これは、重要な性質である。ノイズに対する耐性を容易に評価でき、一般的に、複雑な非線形モデルよりもノイズに対して強い傾向にある。

しかし、このアフィンカメラモデルはあくまでも現実のカメラの近似である。一般には、現実のカメラを透視変換カメラモデルでモデル化するのが普通である。もちろん、レンズの歪み等を考えると透視変換カメラであっても十分でないことも考えられ得るものの、透視変換モデルを正しいカメラモデルであるとしておく。

このアフィンカメラモデルと透視変換カメラモデルの間には大きな違いがある。線形であるか、非線形であるかという違いである。もし、透視変換カメラで撮影された画像から、アフィンカメラで撮影した画像を容易に合成できるならば、3次元復元の操作は線形になり、よい性質を持つと考えられる。また、このような操作は、アフィンカメラと透視変換カメラの中間的な性質を持つとも考えられる。

本論文では、複数の画像が得られた場合に、それを仮想的なアフィンカメラで撮影したと同等な画像を生成し、その生成した画像から、線形な操作で三

次元位置情報の復元を試みる。

2 従来の研究

三次元空間に存在する点(対象点)をカメラで撮影し、対象点の三次元位置を推定する手法は、コンピュータビジョンの中心的な問題である。

これを解決する最も基本的な手法として、立体視があげられる。あらかじめ相互の位置、姿勢のわかっている2台のカメラで対象点を撮影し、その投影像から三角測量の原理で三次元位置を決定するものである。

立体視では、カメラの位置、姿勢、焦点距離などを正確に測定する必要がある。これは、カメラキャリブレーションと呼ばれ、コンピュータビジョンやロボティクスの分野で従来から研究されてきた。このとき、三次元空間と画像の関係の記述法として透視変換による関係が一般に採用されてきた。この透視変換モデルは、一般的なカメラの理想的なモデルであるといえる。しかし、この投影モデルは正確である反面、非線形という性質を持っている。立体視を用いた三次元位置推定が計算誤差や投影点の計測位置誤差に弱いという性質は、この非線形性が原因である。

この透視変換モデルを、性質のよいカメラモデルで近似するという研究も盛んに行われている。カメラモデルの近似をカメラキャリブレーションに応用した研究の初期のものとして、Grembanらの研究[1]が挙げられる。その後、Quanによるアフィンカメラモデルの研究[2]が進んだ。アフィンカメラモデルでは、三次元空間と画像を線形な関係で記述するものである。線形であるため、非線形性に由来する欠点は解消される。一般に、対象自体の厚みが、カメラと対象の間の距離に比べて十分小さい場合は、透視変換モデルの十分良い近似であることが知られている。

このアフィンカメラモデルの三次元位置推定における応用例として、Christy[3]らの研究があげられる。まず、アフィンカメラモデルを用いて、対象

点の三次元位置の近似的な推定を行い、この近似値を初期値として、透視変換カメラから得られる非線形方程式の最適化を行い、対象点の三次元位置をより正確に推定するという手法を提案している。しかし、この方法では非線形方程式の最適化という操作が必要で、簡潔な解法であるとはいいいにくい。

立体視のもうひとつの流れとして、カメラを複数（2台以上）用いる、多眼立体視がある。これは、2台だけを用いる場合に比べ、情報量が増え、より安定に三次元位置を推定できると期待できる。その例としてTomasiらの研究[4]が挙げられる。Tomasiらによると、複数の正射影カメラを仮定できるなら、因子分解法と呼ばれる簡潔な方法で三次元位置の推定が可能である。もし、カメラが正射影カメラ（正確にはアフィンカメラ）でない場合の三次元位置推定は、エピソード拘束と呼ばれる、非線形な拘束条件を課す必要がある。そのため、容易に三次元位置を推定できるというわけではない。

本論文では、複数のカメラで撮影した画像から、仮想的なアフィンカメラで撮影した画像を生成し、三次元位置を線形演算で推定する手法を提案する。

3 カメラモデル

まず、カメラのモデル化について簡単にまとめる。

3.1 カメラモデルの定義

三次元空間の点を P とし、その投影像の画像上での位置を p とする。これらを齊次座標で、それぞれ \mathbf{P}, \mathbf{p} と表現すると、

$$\mathbf{p} \simeq M\mathbf{P}$$

という関係がある。 \simeq は定数倍を除いて等しいことをあらわす。 M は 3×4 の行列で、射影行列と呼ばれ、投影モデル、カメラの位置・姿勢、光学的な性質などをあらわす行列である。

M が次の様に分解できる場合、このようなカメ

ラモデルを透視カメラモデルと呼ぶ。

$$M = \begin{pmatrix} f_u & 0 & u_0 \\ 0 & f_v & v_0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 1 \end{pmatrix}$$

ただし、 f_u, f_v はカメラの焦点距離であり、 u_0, v_0 は画像中心、 \mathbf{R}, \mathbf{t} はカメラの姿勢と位置を表す回転行列と並進ベクトルである。これは、一般的なカメラの理想的なモデルである。

これに対して、 M の第4行第1列から第3列までが0である場合をアフィンカメラモデルと呼ぶ。

$$M_A = \begin{pmatrix} * & * & * & * \\ * & * & * & * \\ 0 & 0 & 0 & * \end{pmatrix}$$

齊次座標表現から一般に用いられているユークリッド座標表現にもどすと、アフィンカメラモデルでは、点 P の三次元での座標を \mathbf{P} とするならば、その投影点 p の画像上での座標 \mathbf{p}_A は、

$$\mathbf{p}_A = C_A \mathbf{P} + \mathbf{d}$$

というかたちで関係付けることが可能である。ただし、 C_A は 2×3 の行列である。ここで、任意の三次元空間の点 P_G とその画像への投影点 p_G を仮定する。 \mathbf{P} と \mathbf{p}_A を、それぞれ \mathbf{P}_G と \mathbf{p}_G を原点するような新しい座標系で表現する。つまり、 $\mathbf{P}' = \mathbf{P} - \mathbf{P}_G$, $\mathbf{p}' = \mathbf{p} - \mathbf{p}_G$ とする。すると、 \mathbf{d} の項が消え、

$$\mathbf{p}' = C_A \mathbf{P}' \quad (1)$$

という簡潔な表現となる。もし、あらかじめ $\mathbf{P}_G, \mathbf{p}_G$ を引いた座標を使うと、三次元空間 \mathbf{P}' と画像 \mathbf{p}' は 2×3 の行列 C_A で関係付けることができることがわかった。この C_A をアフィンカメラ行列と呼ぶことにする。 P_G の例として、もし、多数の点が存在する場合、それらの点の重心を P_G として用いるのが簡便であろう。アフィンカメラの仮定のもとでは、三次元空間での重心 P_G の投影点は、画像上の投影点の重心 p_G と一致する。

以降、あらかじめ、この重心を差し引いた座標系を採用する。

3.2 拡張アフィンカメラ

さて次に、三次元空間と、複数のカメラで得られた複数の画像とのあいだの関係を調べる。

複数 (N 台) のカメラから、三次元空間の点 P を撮影しているとする (図 1)。各カメラでの P の投影点を p_1, \dots, p_N とする。これは、点 P を N 台のカメラで撮影して得られるものであるが、三次元空間の点を $2N$ 次元の「画像」に投影したものと考えることもできる (図 2)。三次元空間を $2N$ 次元の画像に投影するカメラのことを拡張カメラと呼ぶことにしよう。注意しなくてはならないのは、拡張カメラでは、三次元からの次元の縮退は起こらない、冗長な観測系になっていることである。

任意の点 (一般に重心) とその投影像を原点とする座標系を採用する。点 P の座標を \mathbf{P} 、画像 i における点 P の投影点の座標を p_i とし、まとめて、

$$\mathbf{p}^* = \begin{bmatrix} p_1 \\ \vdots \\ p_N \end{bmatrix} \quad (2)$$

と表記することにする。これは、拡張カメラにおける投影点の座標といえる。

もし、各カメラがアフィンカメラであると仮定するならば、カメラ i でのアフィンカメラ行列を C_i とすると、各カメラでの座標 \mathbf{p}_{A_i} は、(1) より、

$$\mathbf{p}_{A_i} = C_i \mathbf{P}$$

という関係がある。全てのカメラについてまとめると、

$$\mathbf{p}^*_{A} = C_V \mathbf{P} \quad (3)$$

と記述できる。ただし、

$$C_V = \begin{bmatrix} C_1 \\ \vdots \\ C_N \end{bmatrix} \quad (4)$$

である。 \mathbf{p}^*_{A} は $2N$ -ベクトルであり、 C_V は $2N \times 3$ の行列である。 C_V は各カメラをアフィンカメラであると仮定した場合の、拡張カメラのアフィンカメラ行列に相当している。

さて、(3) という関係のもとで、仮に C_V が既知であるならば、画像 \mathbf{p}^*_{A} から点 P の三次元位置 \mathbf{P} を解くことができる。それには、 C_V の疑似逆行列を使って、

$$\mathbf{P} = C_V^{\dagger} \mathbf{p}^*_{A} \quad (5)$$

とすればよい。

観測量 \mathbf{p}^*_{A} にノイズが乗っていないのならば、つまり、 \mathbf{p}^*_{A} が正確に観測されている状況では、(5) によって得られる \mathbf{P} は、正確な三次元位置となっている。つまり、この \mathbf{P} は近似的な解ではなく、アフィンカメラの仮定のもとでは、正しい三次元位置の復元が可能である。 \mathbf{p}^*_{A} にノイズが乗っている場合は、(5) の解法は、 $\|\mathbf{p}^*_{A} - C_V \mathbf{P}\|$ を最小に、つまり、拡張アフィンカメラでの画像上の誤差を最小にするような推定となっている。拡張アフィンカメラを仮定するなら、三次元空間の座標 \mathbf{P} と画像 \mathbf{p}^*_{A} の間は線形の関係となっている。線形であるゆえ、解析やノイズに対する耐性などが非線形の場合とくらべて格段に良いことは明らかである。

3.3 透視カメラからアフィンカメラへ

拡張アフィンカメラを仮定できる場合は、その画像から、三次元空間の点の三次元位置を線形なかたちで復元できることがわかった。これは、アフィンカメラを仮定できる場合だけで、透視カメラの場合は、例えノイズが乗っていない場合でも、(5) の方法では、三次元位置を正確には復元できないことは容易に推察できる。しかし、もし、透視カメラの画像から、なんらかの方法を使って、アフィンカメラで撮影したのと等価な画像を生成することができたならば、その画像を使って、三次元位置を線形に推定することが可能なはずである。

3.4 アフィンカメラへの補正

ここでは、透視カメラで得られた画像を、アフィンカメラで撮影した画像と等価な画像に補正する方法を提案する。エピソード幾何の拘束を使うことも考えられるが、ここでは、できるだけ単純な補正方

法を提案する。

以下では、拡張カメラを用いて説明する。

三次元空間の点 P を現実の N 台のカメラで撮影した画像を \mathbf{p}^* とする。また、行列 C_V で与えられる仮想的な拡張アフィンカメラを考え、これで点 P を撮影した画像を \mathbf{p}^*_A とする。 \mathbf{p}^* が \mathbf{p}^*_A と一致するように、 \mathbf{p}^* を補正する(図3)。さまざまな補正方法が考えられるが、ここでは、 \mathbf{p}^*_A の k 番目の要素 $(\mathbf{p}^*_A)_k, (1 \leq k \leq 2N)$ を \mathbf{p}^* の二次多項式で近似することにする。つまり、

$$(\mathbf{p}^*_A)_k = \tilde{\mathbf{p}}^{*T} T_k \mathbf{p}^* \quad (6)$$

というかたちでモデル化する。ただし、 $\tilde{\mathbf{p}}^* = \begin{bmatrix} \mathbf{p}^* \\ 1 \end{bmatrix}$ であり、 T_k は $(2N+1) \times (2N+1)$ の対称行列である。 T の自由度は、全体で $2N \times (2N+1) \times (N+1)$ であるので、透視カメラモデルを十分記述できると考えられる。

ここで注意すべきは、 T はすべての点について共通のものであるということである。 T は現実のカメラの投影法、カメラの位置、姿勢など、カメラの性質や配置によって決定されるパラメータである。そのため、カメラを移動させたり、焦点距離を変化させない限り、同じ T を使い続けることができる。

4 三次元復元手法

次に、具体的に対象点の三次元位置を決定する手順を述べる。これには、二つの段階が必要である。第一段階は、仮想的な拡張アフィンカメラのパラメータ C_V と補正のためのパラメータ T を決定する校正段階である。そして、第二段階は、 T を用いて現実のカメラで得られた画像を拡張アフィンカメラに投影したものと同等になるように補正し、対象点の三次元位置を線形に推定する推定段階である。

1. 校正段階

三次元位置情報を復元するには、拡張アフィンカメラ行列 C_V と補正用の行列 T を求めておく必要がある。この操作はシステムの校正

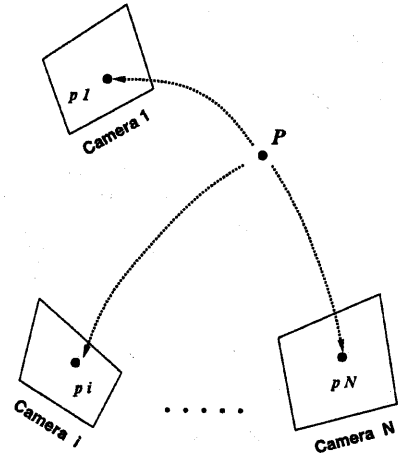


図 1. N 台のカメラで三次元空間の点 P を撮影する。得られるのは、 $2N$ 次元の画像データである。

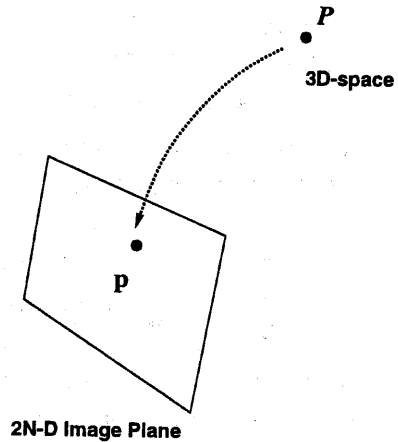


図 2. 拡張カメラの概念図。複数のカメラで撮影した複数の画像を、 $2N$ 次元の一枚の画像であるとみなす。

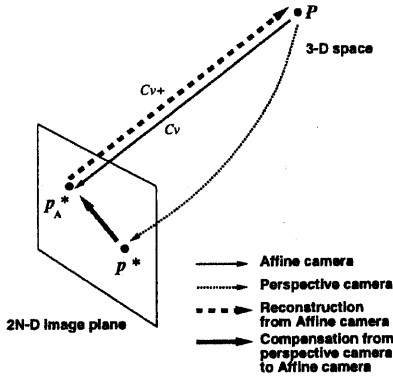


図 3. 透視変換カメラによる投影からアフィンカメラによる投影への補正。 $\mathbf{p}_A^* = \mathbf{p}^{*T} T \bar{\mathbf{p}}^*$ 。アフィンカメラでは、3次元空間から画像、画像から3次元空間は、線形な関係となる。

とみなすことができる。このために、あらかじめ3次元位置のわかっている M 個の参照点 $\mathbf{P}_1, \dots, \mathbf{P}_M$ と、これらの N 台のカメラでの投影像 $\mathbf{p}_1^*, \dots, \mathbf{p}_N^*$ を使う。キャリブレーション用の格子を使っても良いし、任意の点列を使ってもよい。しかし、本手法があくまでも近似的な3次元復元方法であることを考えると、実際に推定したい3次元空間の領域とあまり異なる領域に参照点を置くのがよい。

この参照点列から、まず、仮想的な拡張アフィンカメラ行列 C_V を求める。基本的に C_V はどのような行列でもよい。アフィンカメラで撮影した像は、アフィン変換で、任意のアフィンカメラで撮影した像に変換できるからである。ここでも、あくまでも本手法が近似であるという点を考慮し、実際のカメラの最も良い近似となるように、仮想的なアフィンカメラのパラメータを選ぶのがよい。もし、カメラが全てアフィンカメラであると仮定するならば、

$$C_V[\mathbf{P}_1 \dots \mathbf{P}_M] = [\mathbf{p}_1^* \dots \mathbf{p}_M^*]$$

という関係があるはずである。もちろん、 $\mathbf{P}_1, \dots, \mathbf{P}_M, \mathbf{p}_1^*, \dots, \mathbf{p}_M^*$ は、重心を引いて

正規化したあとの座標である。この方程式を、拡張カメラにおける、全ての点の誤差の総計 $\|[\mathbf{p}_1^* \dots \mathbf{p}_M^*] - C_0[\mathbf{P}_1, \dots, \mathbf{P}_M]\|$ が最小になるように解き、

$$C_V = [\mathbf{p}_1^* \dots \mathbf{p}_M^*]^+ [\mathbf{P}_1, \dots, \mathbf{P}_M] \quad (7)$$

を得る。

次に、 T を求める。 C_V という拡張アフィンカメラで点 P_j が撮影されたとすると、画像上、 $\mathbf{p}_{A_j}^* = C_V \mathbf{P}_j$ に投影される。この $\mathbf{p}_{A_j}^*$ を得るように、 $\bar{\mathbf{p}}_j^*$ を補正する。補正方法は、先に述べたように、 $(\mathbf{p}_{A_j}^*)_k = \bar{\mathbf{p}}_j^{*T} T_k \bar{\mathbf{p}}_j^*$ である。いま、 $\bar{\mathbf{p}}_j^*, \mathbf{p}_{A_j}^*, (1 \leq j \leq M)$ が得られているので、 $T_k, (1 \leq k \leq 2N)$ を求めるには、各 T_k の $(2N+1) \times (N+1)$ 個の要素に関して線形方程式を解くだけで良い。

これで、 M 個の参照点の3次元位置 \mathbf{P}_j と、それらの N 台のカメラでの投影像 \mathbf{p}_{ij}^* から、仮想的なアフィンカメラ行列 C_V 、それに対応する補正行列 $T_k, (1 \leq k \leq 2N)$ が決定された。

2. 3次元位置推定

3次元空間中に存在する点 Q の3次元位置を推定する。校正段階と同じ N 台のカメラでの Q の投影像 \mathbf{q}^* が得られているとする。これから、仮想的な拡張アフィンカメラ C_V で撮影されたのと同等の投影像 \mathbf{q}_A^* を得るように、

$$(\mathbf{q}_A^*)_k = \bar{\mathbf{q}}^{*T} T_k \bar{\mathbf{q}}^*, (1 \leq k \leq 2N)$$

と補正する。そして、この \mathbf{q}_A^* を用いて、

$$\mathbf{Q} = C_V^+ \mathbf{q}_A^* \quad (8)$$

と Q の3次元位置 \mathbf{Q} が推定される。

5 本手法のまとめ

仮想的なアフィンカメラを設定し、そのアフィンカメラで撮影したのと等価な画像を生成することで、3次元位置を線形に推定する方法を提案した。

その手順を以下にまとめる。

1. 校正段階

- (a) 三次元位置のわかっている M 個の参照点 P_1, \dots, P_M を N 台のカメラで撮影する。
- (b) 参照点の重心 P_G 、各カメラでの重心 p^*_G を考慮し、座標を正規化する。
- (c) 参照点 P_1, \dots, P_M の座標 P_1, \dots, P_M と、その投影点の座標 p^*_1, \dots, p^*_M から、仮想的なアフィンカメラ行列 C_V を計算する。
- (d) この仮想的なアフィンカメラで P_1, \dots, P_M を投影し、その投影点の座標を $p^*_{A1}, \dots, p^*_{AM}$ とする。
- (e) 投影点の座標 p^*_i を p^*_{Ai} ($1 \leq j \leq M$) に補正するような T_k ($1 \leq k \leq 2N$) を計算する。

2. 推定段階

- (a) 対象点 Q を、校正段階と同じ N 台のカメラで撮影する。
- (b) 校正段階と同じ P_G, p^*_G を使って、座標を正規化する。
- (c) 校正段階と同じ仮想的なアフィンカメラでの投影点の座標 q^*_A を計算する。
- (d) C_V の一般化逆行列 C^+_V を計算し、 C^+_V と q^*_A から、対象点 Q の座標を、 $Q = C^+_V q^*_A$ と推定する。

6 シミュレーション実験

簡単なシミュレーション実験の結果を示す。

三次元空間中に 64 点の参照点を配置する。カメラは 4 台、これらの参照点群が各カメラの視野に入るように設置した。64 点の参照点を使ってカメラのキャリブレーション (C_V, T を推定) したあと、これら参照点群とは別の、そして、参照点群よ

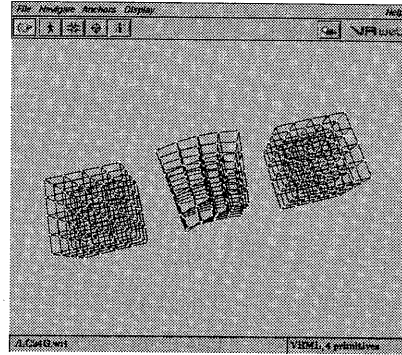


図 4. 4 台のカメラによる三次元位置の推定結果。左：もとの三次元形状。中：アフィンカメラ補正しない場合の三次元形状。右：アフィンカメラ補正をした三次元形状。（三者が重ならないように、ずらして表示している）

りも広い範囲に格子状に配置された 125 点を使って、三次元位置推定を行った。図 4 にその結果を示す。左の点群は、もとの、格子状に配置された対象点である。これを、透視投影カメラで撮影して得られた投影像を、アフィンカメラ補正せずに、そのまま、(8) を使って三次元位置を推定したのが真ん中の点群である。これは、透視投影カメラをアフィンカメラで近似したことになる。いくらか歪んだ形状が再現されているのがわかる。右の点群は、アフィンカメラ補正を施して、(8) により三次元位置を推定したものである。左の点群と形状が一致することが見てとれる。

7 おわりに

仮想的なアフィンカメラを設定し、そのアフィンカメラで撮影したのと等価な画像を生成することで、三次元位置を線形に推定する方法を提案した。この操作は、カメラモデルに対して変更を加えたことと等価であり、透視変換カメラとアフィンカメラの両方の特徴をうまく融合していると考えられる。

今後の問題点としては、透視変換カメラモデル

を、多項式(2次式)で近似することが、果たしてどのくらい正確であるのかという点である。1次式ではアフィンカメラと同等。最も良いのは透視変換カメラの拘束を直接用いることであるが、これでは複雑になる。その妥協点として2次式の形を採用した。ところが、カメラの数が多くなると、この2次式の自由パラメータの数が多くなる。つまり、このモデルの自由度が大きくなるということである。透視変換カメラのモデル化に本来必要な自由度以上に自由度がある状況では、キャリブレーション時にノイズが混入してしまった場合、このノイズを含めたシステムを大きな自由度を利用してモデル化してしまうという現象が起きる可能性がある。ノイズの大きな影響を受ける可能性がある。そのような場合は、線形近似以上、2次式による近似以下の自由度を持つ近似方法を提案する必要があるだろう。

いずれにしても、三次元空間の位置を非線形な形で直接計算するのではなく、いったん、画像上で非線形な補正をすることで、三次元復元に関しては線形に推定するという点が、従来と異なる特徴である。つまり、その間接的な計測量である画像を操作することで、三次元空間の位置を操作していることになっている。

参考文献

- [1] Keith D.Gremban, Charles E. Thorpe, and Takeo Kanade. Geometric camera calibration using systems of linear equations. In *IEEE International Conference on Robotics and Automation*, pages 562-567, 1988.
- [2] Long Quan. Self-calibration of an affine camera from multiple views. *International Journal of Computer Vision*, 19(1):93-105, 1996.
- [3] S.Christy and R.Horand. Euclidean shape and motion from multiple perspective views by affine iterations. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 18(11):1098-1104, 1996.
- [4] Tomasi and Kanade. Shape and motion from image streames under orthography: A factorization method. *International Journal of Computer Vision*, 9(2):137-154, 1992.