

頑健なヤコビ核主成分分析に向けて

藤木 淳[†] 赤穂 昭太郎[†]
日野 英逸^{††} 村田 昇^{††}

元来の核主成分分析では特徴空間におけるユークリッド計量に基づいた最小二乗推定によって主たる部分空間を求めるが、ヤコビ核主成分分析では入力空間におけるユークリッド計量に基づいて計算される重みを用いた重み付き最小二乗推定によって主たる部分空間を求める。一般に（重み付き）最小二乗推定量は外乱値に弱いので、本稿では RANSAC と χ^2 検定を利用して外乱値を検出することによりヤコビ核主成分分析の頑健化を行なう。

Towards robust Jacobian kernel PCA

JUN FUJIKI,[†] SHOTARO AKAHO,[†] HIDEITSU HINO^{††}
and NOBORU MURATA^{††}

Conventional kernel principal component analysis finds the principal subspace of the data by means of least squares estimation with respect to the metric defined in the feature space. On the other hand, Jacobian kernel principle component analysis (JKPCA) finds the principal subspace based on the metric in the input space by weighted least squares associated with the Jacobian of the feature map. For some applications such as image processing, it is more natural to utilize the metric in the input space, however, the weighted least squares estimation is sometimes too sensitive to outliers. To overcome this drawback of JKPCA and make it robust to outliers, an outlier detection method based on random sampling consensus (RANSAC) and χ^2 test is discussed, and its validity is confirmed by numerical experiments.

1. はじめに

観測データの構造の把握はデータ解析における基本的かつ重要な手順である。データの主構造の最も基本的な決定手法は主成分分析 (principal component analysis ; PCA) である。PCA は、データをユークリッド (欧幾里得 ; 欧氏) 空間の点とみなし、データの真値は欧氏空間内のアフィン空間上に分布するという仮定の下で有効な手法である。そして PCA ではデータと主構造とのずれである誤差を欧氏距離で測定し、誤差の二乗和を最小とする最小二乗 (least squares ; LS) 基準を用いるが、これは誤差が等方正規分布に従うという仮定の下では最尤推定を与える。

しかしデータの真値や主構造がアフィン空間上に分布しない場合には主構造としてアフィン空間を仮定する PCA は不向きである。そこで非線型主成分分析 (nonlinear PCA ; NLPCA) が提案されてきた。

NLPCA の考え方は、非線型構造をもつデータをアフィン空間上に分布させるために特徴空間 (feature space) と呼ばれる高次元空間に特徴写像 (feature map) と呼ばれる写像を行ない、その後通常 PCA を行なう所にある。このとき、データの構造は特徴写像によってアフィン空間に写像される空間に限定されるため、データ構造の理解と特徴写像はほぼ等価である。よって特徴写像が理論的に既知の場合は問題がないが、未知の場合には、写像が先か構造が先かという鶏と卵の問題が生じる。しかしデータの構造を十分に近似できる写像が求まれば良いならば、多数の基底で表現される非常に高次元への写像を考えることにより鶏と卵の問題は近似的に解決される。

一般に特徴空間が高次元の場合、特徴空間における高次元ベクトルの計算量が多くなる。この問題を解決したのがカーネルトリック^{1), 2)} である。カーネルトリックとは、NLPCA に特徴写像の表現は不要で、特徴空間での誤差を測るための計量 (内積) を与える対称な正定値カーネル関数があれば十分であるという議論のすりかえのことである。この議論のすりかえにより、計算量が特徴写像の次元数ではなくデータ数に依

[†] 産業技術総合研究所
National Institute of Advanced Industrial Science and Technology
^{††} 早稲田大学 Waseda University

存することとなり非常に高い次元への特徴写像を考えても計算量が破綻しにくくなる。そしてまたガウシアンカーネルのような無限次元空間への写像を扱うこともできる。このカーネルトリックを用いた NLPCA を **核主成分分析 (kernel PCA; KPCA)**^{4),18)} と呼ぶ。

さて、KPCA¹⁸⁾ は特徴空間における欧氏距離の LS 基準によるあてはめ問題である。ゆえにあてはめ結果による入力データの推定値は特徴空間におけるあてはめ結果への最短距離であり、入力空間におけるあてはめ結果への最短距離とは限らない。そこでコンピュータビジョンでは入力空間における距離の LS 基準によるあてはめ問題が研究されてきた。入力空間が欧氏空間の場合^{2),5),14),15),17),20)} や球面の場合⁹⁾ が提案され、そして一般的にリーマン空間の場合を核関数を用いて表現した **ヤコビ核主成分分析 (Jacobian KPCA; JKPCA)**¹⁰⁾ が提案された。

JKPCA は特徴写像を各観測データ点附近で Jacobi 行列を利用してアフィン写像に近似している。このアフィン写像により観測データ点附近における入力空間と特徴空間の計量の対応を近似し、入力空間の誤差を特徴空間の誤差によって近似的に表現している。このことにより、入力空間の LS 推定は特徴空間の重み付き LS 推定により近似される。この近似は **等計量線 (equi-metric curve)**、つまり入力空間においてデータから等距離にある閉曲線の特徴写像による像を超楕円体によって近似していることに相当する。しかし一般に特徴写像は非線型写像であるから、この近似は誤差が十分に小さい範囲内でしか成立せず、**外乱値 (outlier)**、つまり誤差の大きいデータに対しては近似誤差と実際の誤差の乖離が大きく、外乱値を含む場合の推定結果は一般的に悪くなる。また LS 推定自体モデルからデータの方向に過剰に適合し外乱値が推定に悪影響を及ぼすことが知られている。

そこで外乱値の影響を抑えるために M 推定量が提案された¹³⁾。M 推定の基本的な考え方は、誤差分布が正規分布から少々ずれたとしても有効性がそれほど落ちないように外乱値の影響を制御した推定量を作る所にある。例えば平均の推定において誤差がラプラス分布に従うと仮定すると、最尤推定量は標本平均ではなく標本中央値となるというように、直感的には、誤差分布が M 推定量から導かれる“正規分布関数と似た”確率分布 (M 推定量によっては正規化できず確率分布と見做せない場合もある) に従うと仮定した場合の最尤推定として捉えることができる。

本稿で提案する **頑健なヤコビ核主成分分析 (Robust JKPCA; R-JKPCA)** はランザック (ran-

dom sample consensus; RANSAC)⁶⁾ と χ^2 検定の組合せによって外乱値を除去し、許容値のみで推定を行なう¹⁹⁾ 手法であるが、これも一種の M 推定とみなせる (対応する確率分布は存在しない)。

提案手法である R-JKPCA の有効性は、人工データを用いた実験によって評価される。

2. ヤコビ核主成分分析

入力空間である m 次元空間 \mathcal{R} の点列に部分空間をあてはめる際、この点列を特徴空間である n 次元ヒルベルト空間 \mathcal{H} に写像し、 \mathcal{H} における線型あてはめに帰着することが多い。このあてはめを \mathcal{H} における LS 基準で推定するのが KPCA である。しかし入力空間に自然な計量が定義されている場合に特徴空間における計量を基準として推定を行うことは好ましくない³⁾ ため、入力空間の計量における近似 LS 基準で部分空間をあてはめる JKPCA¹⁰⁾ が提案された。

2.1 入力空間の計量と特徴空間の計量

入力空間を m 次元リーマン空間 \mathcal{R} とし、 \mathcal{R} 上の点 \mathbf{x} におけるリーマン計量を $G_{\mathbf{x}}$ とする。また、観測データ点 $\{\mathbf{x}_{[d]}\}_{d=1}^D$ とする。JKPCA では観測データ点 $\mathbf{x}_{[d]}$ 附近の空間を計量が $G_{\mathbf{x}_{[d]}}$ で一定であるアフィン空間で近似する。つまり $\mathbf{x} = \mathbf{x}_{[d]} + \delta\mathbf{x}$ なる点 \mathbf{x} と観測データ点 $\mathbf{x}_{[d]}$ の距離 r^p を $(r^p)^2 = (\delta\mathbf{x})^T G_{\mathbf{x}_{[d]}} (\delta\mathbf{x})$ で近似する。点 \mathbf{x} を n 次元ヒルベルト空間である特徴空間 \mathcal{H} に特徴写像 $\phi: \mathbf{x} \mapsto \phi(\mathbf{x})$ で射影する。ここで特徴写像 ϕ のヤコビ行列 J_{ϕ} を用いると、特徴写像は観測データ点附近で

$$\mathbf{x} \mapsto \phi_{[d]} + J_{\phi_{[d]}} (\mathbf{x} - \mathbf{x}_{[d]})$$

($\phi_{[d]} = \phi(\mathbf{x}_{[d]})$) と近似でき、

$$(r^p)^2 = (\delta\phi(\mathbf{x}))^T G_{\phi_{[d]}} \delta\phi(\mathbf{x})$$

$$\begin{bmatrix} \phi(\mathbf{x} + \delta\mathbf{x}) = \phi(\mathbf{x}) + \delta\phi \\ G_{\phi_{[d]}} = (J_{\phi_{[d]}}^+)^T G_{\mathbf{x}_{[d]}} J_{\phi_{[d]}}^+ \\ G_{\phi_{[d]}}^{-1} = J_{\phi_{[d]}} G_{\mathbf{x}_{[d]}}^{-1} J_{\phi_{[d]}}^T \end{bmatrix}$$

(X^+ は X のムーア・ペンローズ逆行列) と近似できる。

2.2 特徴空間における線型あてはめ

特徴空間のデータ $\phi_{[d]} = \phi(\mathbf{x}_{[d]})$ に対して $n-1$ 次元アフィン空間 $\mathbf{a}^T \phi + b = 0$ をあてはめるには、

$$\mathcal{E}(\mathbf{a}) = \sum_{d=1}^D \frac{(\mathbf{a}^T \phi_{[d]} + b)^2}{\mathbf{a}^T G_{\phi_{[d]}}^{-1} \mathbf{a}} \quad (1)$$

を最小にする \mathbf{a} , b を求めれば良い¹⁰⁾。ここでは

$$\tilde{\phi} = \begin{pmatrix} \phi \\ 1 \end{pmatrix}, \quad \tilde{\mathbf{a}} = \begin{pmatrix} \mathbf{a} \\ b \end{pmatrix}, \quad \tilde{G}_{\phi_{[d]}} = \begin{pmatrix} G_{\phi_{[d]}} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{pmatrix}$$

と置いて $n+1$ 次元空間の原点を通る n 次元線型部分空間 $\tilde{\mathbf{a}}^\top \tilde{\phi} = 0$ の推定問題に帰着させる, つまり

$$\mathcal{E}(\tilde{\mathbf{a}}) = \sum_{d=1}^D \frac{\tilde{\mathbf{a}}^\top \left[\tilde{\phi}_{[d]} \tilde{\phi}_{[d]}^\top \right] \tilde{\mathbf{a}}}{\tilde{\mathbf{a}}^\top \tilde{G}_{\phi_{[d]}}^+ \tilde{\mathbf{a}}} \quad (2)$$

の最小にする $\tilde{\mathbf{a}}$ を求めれば良く^{2),14),20)}, 式 (2) の最小化を核関数を用いて行なう.

2.3 核関数による表現

本稿では, 核関数として微分可能なものを考え, 核関数 $k(\mathbf{x}, \mathbf{y}) = \tilde{\phi}(\mathbf{x})^\top \tilde{\phi}(\mathbf{y})$ だけでなく, その微分

$$\mathbf{k}(\mathbf{x}, \mathbf{y}) = \frac{\partial k(\mathbf{x}, \mathbf{y})}{\partial \mathbf{x}^\top} = J_{\tilde{\phi}}(\mathbf{x})^\top \tilde{\phi}(\mathbf{y})$$

も用いる. これをヤコビ核 (Jacobian kernel) と呼ぶ. また \mathbf{a} の存在範囲として $\tilde{\phi}_{[d]}$ の線型結合

$$\mathbf{a} = \sum_p \alpha_{[d]} \tilde{\phi}_{[d]} = \tilde{\Phi} \boldsymbol{\alpha}$$

だけを考える. ここで $D \times D$ 行列 \mathcal{K} を

$$(\mathcal{K})_{ij} = k(\mathbf{x}_{[i]}, \mathbf{x}_{[j]}), \\ \mathcal{K} = \begin{pmatrix} \mathcal{K}_{[1]} & \cdots & \mathcal{K}_{[D]} \end{pmatrix}$$

で定義し, また, $D \times m$ 行列 $\mathcal{K}_{[d]}$ を

$$\mathbf{k}_{[i][j]} = \mathbf{k}(\mathbf{x}_{[i]}, \mathbf{x}_{[j]}), \\ \mathcal{K}_{[d]} = \begin{pmatrix} \mathbf{k}_{[d][1]} & \cdots & \mathbf{k}_{[d][D]} \end{pmatrix}^\top$$

で定義すると,

$$\mathcal{K}_{[d]} = \tilde{\Phi}^\top \tilde{\phi}_{[d]}, \quad \mathcal{K}_{[d]} = \tilde{\Phi}^\top J_{\tilde{\phi}_{[d]}}$$

が成立し, このとき式 (2) は, 写像 ϕ を含まない

$$\mathcal{E}'(\boldsymbol{\alpha}) = \sum_{d=1}^D \frac{\boldsymbol{\alpha}^\top \mathcal{K}_{[d]} \mathcal{K}_{[d]}^\top \boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top \mathcal{K}_{[d]} \tilde{G}_{\mathbf{x}_{[d]}}^+ \mathcal{K}_{[d]}^\top \boldsymbol{\alpha}} \quad (3)$$

となり, これを最小にする $\boldsymbol{\alpha}$ を求めれば良い. ここで

$$\tilde{G}_{\mathbf{x}_{[d]}} = \begin{pmatrix} G_{\mathbf{x}_{[d]}} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{pmatrix}$$

である.

2.4 解法とアルゴリズム

式 (3) の最小化アルゴリズム²⁾ を紹介する. $\boldsymbol{\alpha}$ の近似値 $\hat{\boldsymbol{\alpha}}$ が得られたとき,

$$\mu_{[d]} = \hat{\boldsymbol{\alpha}}^\top \mathcal{K}_{[d]} \tilde{G}_{\mathbf{x}_{[d]}}^+ \mathcal{K}_{[d]}^\top \hat{\boldsymbol{\alpha}}, \\ \Lambda = \text{diag}\{\mu_{[1]}, \dots, \mu_{[D]}\}$$

とすると $\mathcal{E}'(\boldsymbol{\alpha}) \approx \boldsymbol{\alpha}^\top \mathcal{K} \Lambda^{-1} \mathcal{K} \boldsymbol{\alpha}$ と近似できるので,

$$\boldsymbol{\alpha} = \text{UnitMinEigVec}[\mathcal{K} \Lambda^{-1} \mathcal{K}]$$

(UnitMinEigVec[X] は正方行列 X の最小固有値に対応する単位固有ベクトル) となる.

以下のアルゴリズムにおいて右肩の $[k]$ によって k ステップ目の値を表すものとする. 初期値は右肩が $[0]$ となる (初期値の設定は後述).

- (1) 初期値 $\mu_{[d]}^{[0]}$ ($d = 1, \dots, D$) の設定
- (2) 収束するまで (a), (b) を繰り返す
 - (a) $\hat{\boldsymbol{\alpha}}^{[k+1]} = \text{UnitMinEigVec}[\mathcal{K}(\Lambda^{[k]})^{-1} \mathcal{K}]$
 - (b) $\mu_{[d]}^{[k+1]} = (\hat{\boldsymbol{\alpha}}^{[k+1]})^\top \mathcal{K}_{[d]} G_{\mathbf{x}_{[d]}}^+ \mathcal{K}_{[d]}^\top (\hat{\boldsymbol{\alpha}}^{[k+1]})$ で $\{\mu_{[d]}^{[k+1]}\}_{d=1}^D$ を更新.

2.5 初期値

上記アルゴリズムの初期値の設定方法として以下のような方法がある.

2.5.1 核主成分分析 (KPCA)

$\{\mu_{[d]}^{[0]} = 1\}_{d=1}^D$ の場合²⁾ で

$$\mathcal{E}'(\boldsymbol{\alpha}) \approx \boldsymbol{\alpha}^\top \mathcal{K}^2 \boldsymbol{\alpha}$$

を最小化する $\hat{\boldsymbol{\alpha}} = \text{UnitMinEigVec}[\mathcal{K}^2]$ を初期値とする.

2.5.2 欧氏化 (Euclideanization)

欧氏化^{7),10),11)} とは特徴写像によるデータ点附近の線分長の平均拡大率の逆数をデータ点に重みとして与える手法であり,

$$\mathcal{E}'(\boldsymbol{\alpha}) \approx \boldsymbol{\alpha}^\top \left(\mathcal{K} D^{\frac{1}{m}} \mathcal{K} \right) \boldsymbol{\alpha},$$

$$D = \text{diag} \left\{ \text{Det } G_{\phi_{[1]}}, \dots, \text{Det } G_{\phi_{[D]}} \right\}$$

$$\text{Det } G_{\phi_{[d]}} = \det \left\{ J_{\phi_{[d]}}^+ (J_{\phi_{[d]}}^+)^T \right\} \cdot \det G_{\mathbf{x}_{[d]}}$$

(Det X で X の 0 でない固有値の積を表す) を最小化する $\hat{\boldsymbol{\alpha}} = \text{UnitMinEigVec}[\mathcal{K} D^{\frac{1}{m}} \mathcal{K}]$ を初期値とする.

2.5.3 Taubin 法²⁰⁾

各データ点で異なる計量を核空間で平均した

$E[G_{[d]}^{\text{kernel}}] = \sum_{d=1}^D \mathcal{K}_{[d]} G_{\mathbf{x}_{[d]}}^+ \mathcal{K}_{[d]}^\top$ で近似する手法であり,

$$\mathcal{E}'(\boldsymbol{\alpha}) \approx \frac{\boldsymbol{\alpha}^\top \mathcal{K}^2 \boldsymbol{\alpha}}{\boldsymbol{\alpha}^\top \left(E[G_{[d]}^{\text{kernel}}] \right) \boldsymbol{\alpha}}$$

を最小化する $\hat{\boldsymbol{\alpha}}$, つまり一般化固有値問題を

$$\mathcal{K}^2 \boldsymbol{\alpha} = \lambda \left(E[G_{[d]}^{\text{kernel}}] \right) \boldsymbol{\alpha}$$

の最小固有値に対応する固有ベクトルを初期値とする.

2.6 ヤコビ核主成分分析

KPCA は $\mathcal{K}^2 = \mathcal{K} \mathcal{K}$ の固有ベクトルを固有値の大きい順に並べて新しい座標軸を作る手法であり, JKPCA

は $\mathcal{K}(\Lambda^{[\infty]})^{-1}\mathcal{K}$ の固有ベクトルを固有値の大きい順に並べて新しい座標軸を作る手法である。よって特徴ベクトル $\tilde{\phi}_{[d]}$ に対して重み

$$(\mu_{[d]}^{[\infty]})^{-\frac{1}{2}} = \|\mathcal{K}_{[d]}^T \alpha\|_{G_{[d]}^+}$$

を与えたときの重みつき PCA である。この重みはデータ点 $\mathbf{x}_{[d]}$ 附近の計量と $\tilde{\phi}_{[d]}$ 附近の計量を比べ、 $\tilde{\phi}_{[d]}$ 附近に局所的に入力空間の計量を反映させたものである。

3. ヤコビ核主成分分析の頑健化

3.1 入力空間の近似誤差

JKPCA は入力空間における誤差 $r_{[d]}$ を

$$r_{[d]}^2 \approx \frac{\alpha^T \mathcal{K}_{[d]} \mathcal{K}_{[d]}^T \alpha}{\alpha^T \mathcal{K}_{[d]} \tilde{G}_{[d]}^+ \mathcal{K}_{[d]}^T \alpha} \quad (4)$$

(α は JKPCA の収束時の値) と近似しており、この近似誤差

$$R_{[d]} = \sqrt{\frac{\alpha^T \mathcal{K}_{[d]} \mathcal{K}_{[d]}^T \alpha}{\alpha^T \mathcal{K}_{[d]} \tilde{G}_{[d]}^+ \mathcal{K}_{[d]}^T \alpha}} \quad (5)$$

を判断基準として JKPCA を頑健化する。

3.2 RANSAC

RANSAC⁶⁾ は、ランダムサンプリングに基づくロバストなモデル作成手法であり、

- (1) 部分空間を形成する最小のデータ数を用いて部分空間をあてはめる。
- (2) あてはめ結果に対して、全データの各々が許容値かどうかを検定し、許容値の個数を数える。
- (3) ランダムランプリングを一定回数繰り返し、許容値の個数が最大となるものを選ぶ。
- (4) 選ばれた個数が最大となる部分空間を用いて全データの各々が例外値かどうかを検定し、例外値を取り除く。
- (5) 例外値を除いたデータ全てから部分空間をあてはめる。

という手続きによって部分空間をあてはめる手法である。

データ量が十分にあるとし、汚染率 (外乱値の割合) が α 、部分空間を作成するのに必要なデータ数を m 個であるとする。このとき m 個のデータが全て許容値である確率が β 以上となるためのランダムサンプリングの回数 k は

$$k \geq \frac{\log_{10}(1 - \beta)}{\log_{10}\{1 - (1 - \alpha)^m\}} \quad (6)$$

をみます。

3.3 MAD 推定量

中央絶対偏差 (median absolute deviation ; MAD) とは、1 次元実データ $\{x_n\}_{n=1}^N$ に対して

$$\text{mad}[x_n] = \text{med} [|x_n - \text{med}(x_n)|]$$

によって定義される統計量である。標準偏差の定義

$$\sigma = \sqrt{\text{Var}[x_n]} = \sqrt{E[(x_n - E[x_n])^2]}$$

と比較してみればわかるように、MAD は標準偏差と同様、確率分布の尺度助変数となる。

ここで $x_n \sim N(0, \sigma^2)$ のとき、 $N \rightarrow \infty$ における漸近的な振舞いを考えると、

$$E[x_n] \simeq 0, \quad \text{med}[x_n] \simeq 0$$

が成立するので、

$$\text{mad}[x_n] \simeq \text{med} [|x_n|], \quad \sigma^2 \simeq E[x_n^2]$$

となる。一方

$$\frac{x_n^2}{\sigma^2} \sim \chi^2(1)$$

であるから、

$$\text{med}[x_n^2] \simeq \chi_{0.5}^2(1) \cdot \sigma^2$$

となり、

$$\sigma \simeq \frac{1}{\sqrt{\chi_{0.5}^2(1)}} \text{mad}[x_n] \approx 1.4826 \text{med} [|x_n|]$$

となる¹⁶⁾。ここで $\chi_{\alpha}^2(L)$ は自由度 L の χ^2 分布の 100 α % 点である。

これと同様に、 L 次元データ $\{\mathbf{x}_n\}_{n=1}^N$ が等方正規分布 $\mathbf{x}_n \sim N(\mathbf{0}_L, \sigma^2 \mathbf{I}_L)$ に従うと仮定すると、

$$\frac{\|\mathbf{x}_n\|^2}{\sigma^2} \sim \chi^2(L)$$

であるから、漸近的に

$$\sigma \simeq \frac{1}{\sqrt{\chi_{0.5}^2(L)}} \text{med} [\|\mathbf{x}_n\|] \quad (7)$$

が成立する。

本稿では m 次元入力空間を n 次元特徴空間に射影した後に n 次元特徴空間における k 次元あてはめ問題を考える。このときモデルからのずれである誤差ベクトルは $L = n - k$ 次元空間のベクトルとなる。しかし、入力空間の像は特徴空間内の m 次元以下の図形へと射影されるため、誤差ベクトルの分布には何らかの偏りが生ずる。しかし問題の単純化のため、誤差ベクトルの分布が等方正規分布 $N(\mathbf{0}_L, \sigma^2 \mathbf{I}_L)$ に従うと仮定して、mad 推定量と σ の関係を利用し、誤差ベクトルのノルムの分散を推定する。

3.4 RANSAC と χ^2 検定

RANSAC と χ^2 検定による外乱値除去手法¹⁹⁾ を JKPCA に以下のように適用する。本稿では、検定に用いる誤差分散の推定量 $\hat{\sigma}^2$ の定め方も与える。

3.4.1 閾値となる誤差分散の決定

- (1) RANSAC で部分空間のあてはめて残差の中央値を計算する、という試行を式 (6) で定まる回数行なう、
- (2) 全ての試行のうち中央値が最小となるものに対して式 (7) から誤差分散の推定量 $\hat{\sigma}$ を求める、

3.4.2 Robust JKPCA (R-JKPCA)

- (1) n 次元特徴ベクトルの集合 $\{\Phi(\mathbf{x}_{[d]})\}_{d=1}^D$ からランダムに $\{\Phi(\mathbf{y}_{[m]})\}_{m=1}^M$ を選び JKPCA により $\tilde{\alpha}$ を求め、全データに対して式 (5) の近似残差 $R_{[d]}$ を計算し、 $R_{[d]}^2 < \hat{\sigma}^2$ なる個数を S とする、
- (2) (1) を反復して S を最大とする主成分部分空間を選び、近似残差を有意水準 $100(1 - \gamma)\%$ で検定、つまり $R_{[d]}^2 \geq \chi_{\gamma}^2(L) \cdot \hat{\sigma}^2$ をみたすデータを外乱値として除去する、
- (3) 許容値から JKPCA で主成分部分空間を推定する、

4. 実験

4.1 2次曲線のあてはめ

2次曲線のあてはめ問題を用いて提案手法の有効性を示す。

許容値は放物線 $y = x^2$ 上の点を x 座標が一様分布 $U[-3, 3]$ に従うように生成し、各座標に切断正規分布 (truncated normal distribution) $*$

$$N_{[-2.32\sigma, 2.32\sigma]}(0, \sigma^2)$$

($\sigma = 0.2$) から生成した誤差を添加したもの、外乱値は長方形内の一様分布 $U[-6, 6] \times [-1.5, 10.5]$ から生成したデータ $(x_{[d]}, y_{[d]})$ のうち、 $y = x^2$ との近似距離 $R_{[d]}$ が 3σ 以上のものとした。

汚染率を $\alpha\%$ とし、許容値を $100(1 - \alpha)$ 点、外乱値を 100α 点 (合計 100 点) 生成し、2次曲線のあてはめを行なう。

2次曲線は平面上の 5 点によって決定されるので本稿の実験では $m = 5$ であり、 $\beta = 0.99$ とすると、式 (6) により汚染率 $\alpha = 0.462 \dots$ のときに $k = 100$ となるので、RANSAC におけるランダムサンプル数は 100 とする。また近似誤差は有意水準 1% ($\gamma = 0.99$) で検定し、各汚染率に対して 100 回のあてはめを行

* 平均 μ 、分散 σ^2 の正規分布のうち区間 $[A, B]$ の部分を正規化した切断正規分布を $N_{[A, B]}(\mu, \sigma^2)$ で表すこととする。この切断正規分布の密度関数は $A < x < B$ において

$$\frac{\frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right)}{\Phi\left(\frac{B - \mu}{\sigma}\right) - \Phi\left(\frac{A - \mu}{\sigma}\right)}$$

となる。ここで $\phi(x)$ 及び $\Phi(x)$ は標準正規分布の密度関数及び累積分布関数である。

なって統計的処理を行なった。

4.2 外乱値の検出精度

データが許容値であることを帰無仮説、外乱値であることを対立仮説とする。このとき、許容値を外乱値と判定することを**第一種過誤**、外乱値を許容値と判定することを**第二種過誤**と呼ぶ。

表 1 に汚染率 (外乱値の割合) に対する 2 つの過誤の割合を**平均値 \pm 標準偏差**の形で記載した。表 1 により、汚染率が 40% 以下でなら 9 割以上のデータが正しく分類できている。

表 1 汚染率と過誤の割合 (%) : 100 点

汚染率	第一種過誤	第二種過誤
0%	2.34 \pm 1.66	—
10%	3.73 \pm 3.08	2.80 \pm 5.87
20%	4.94 \pm 4.37	3.40 \pm 5.12
30%	6.19 \pm 5.33	3.73 \pm 4.06
40%	7.78 \pm 6.57	5.60 \pm 5.51
50%	24.82 \pm 11.72	15.10 \pm 8.69

4.3 あてはめ結果

図 1 は汚染率 10% のデータへのあてはめ結果である。図 2 は汚染率 30%、50% の場合の R-JKPCA によるあてはめ結果である。これら図において黒丸は許容値として生成したデータであり、罰点は例外値として生成したデータである。また点線は $y = x^2$ であり、実線はあてはめ結果である。図 1 より JKPCA は外乱値の影響で悪い結果だが、R-JKPCA は外乱値の影響を軽減していることがわかり、図 2 より汚染率が 50% となると計算が破綻することもあるが、概ねうまくあてはまっていることがわかる。

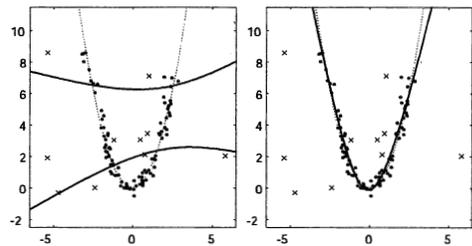


図 1 JKPCA (左) と R-JKPCA (右) (100 点、汚染率 10%)

4.4 許容値に対する近似誤差と真の誤差

あてはめにおいて、検定により許容値と判定されたデータに対し、近似誤差 $R_{[d]}$ 及びデータと $y = x^2$ の距離 (真の距離と呼ぶ) が汚染率に対してどれだけ頑健か調べた。図 3 は近似誤差と真の誤差の平均と分散

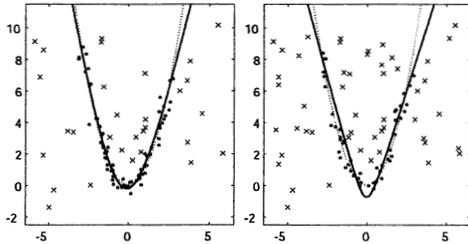


図2 R-JKPCA;汚染率 30% (左), 50% (右)

である。実線、点線はそれぞれ近似誤差及び真の誤差の平均であり、左及び右の縦棒は近似距離(左)及び真の距離(右)それぞれの標準偏差を表す。実験により、汚染率の増加に比べて許容値一個あたりの誤差が抑えられており、提案手法の頑健性が確認できた。

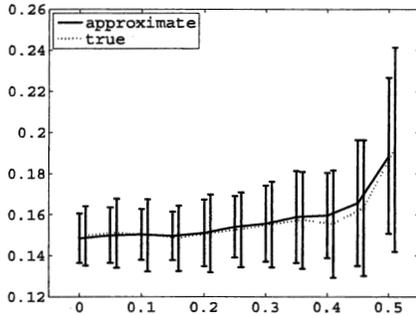


図3 汚染率と許容値あたりの近似誤差

参考文献

- 1) M. Aizerman, É. Braverman and L. Rozonoér, "Theoretical foundations of the potential function method in pattern recognition learning," *Automation and Remote Control*, **25**:821-837, 1964.
- 2) S. Akaho, "Curve fitting that minimizes the mean square of perpendicular distances from sample points," *SPIE, Vision Geometry II*, 1993.
- 3) 赤穂昭太郎, "入力空間でのマージンを最大化するサポートベクターマシン," *信学論 D-II*, **J86-D-II(7)**:934-942, 2003.
- 4) 赤穂昭太郎, "カーネル多変量解析-非線形データ解析の新しい展開-, 岩波書店, 2008.
- 5) W. Chojnacki, M. J. Brooks, A. vanden Hangel and D. Gawley, "On the fitting of surface to data with covariances," *IEEE TPAMI*,

- 22(11)**:1294-1303, 2000.
- 6) M.A.Fischer and R.C.Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Comm. ACM*, **24**:381-395, 1981
- 7) 藤木淳, 赤穂昭太郎, "球面上の点列への小円あてはめ~カメラ運動の平滑化に向けて~, 信学技報 PRMU2004-149:91-96, 2004(12).
- 8) J.Fujiki and S.Akaho, "Spherical PCA with Euclideanization," *Subspace 2007(ACCV07)*.
- 9) 藤木淳, 赤穂昭太郎, "球面最小二乗法による球面上の曲線あてはめ," *Subspace 2008(MIRU2008)*.
- 10) 藤木淳, 赤穂昭太郎, "入力空間での計量に基づいた核主成分分析," *信学技報*, **108(327)**:69-74, 2008.
- 11) 藤木淳, 赤穂昭太郎, 日野英逸, 村田昇, "主成分曲線のあてはめによる放射対称歪曲の較正," *信学技報*, **108(363)**:13-18, 2008.
- 12) R.Hartley and A.Zisserman. *Multiple view geometry in computer vision*. Cambridge University, Cambridge, 2nd edition, 2003.
- 13) P. J. Huber, "Robust estimation of a location parameter," *Ann. Math. Stat.*, **35**:73-101, 1964.
- 14) 金谷健一, 菅谷保之, "幾何学的あてはめの厳密な最尤推定の統一的計算法," *情処研報*, 2008-CVIM-164-3:17-24, 2008.
- 15) Y. Leeden and P. Meer, "Heteroscedastic regression in computer vision: problems with bilinear constraint," *IJCV*, **37(2)**:127-150, 2000.
- 16) R. J. Rousseeuw and A. M. Leroy, *Robust Regression and Outlier Detection*, John Wiley & Sons, NY, 1987.
- 17) P.D.Sampson, "Fitting conic sections to 'very scattered' data: an iterative refinement of the Bookstein algorithm," *Comput. Vision, Graphics, and Image Processing*, **18**:97-108, 1982.
- 18) B.Schölkopf, A.Smola and K.-R.Müller, "Non-linear component analysis as a kernel eigenvalue problem," *Neural Computation*, **10**:1299-1319, 1998.
- 19) Y.Sugaya and K.Kanatani, "Outlier removal for motion tracking by subspace separation," *IEICE Trans. Inf.&Syst.*, **E86-D(6)**:1095-1102, 2003.
- 20) G.Taubin, "Estimation of planar curves, surfaces, and nonplanar space curves defined by implicit equations with applications to edge and range image segmentation," *IEEE TPAMI*, **13(11)**:1115-1138, 1991.
- 21) V.Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, 1995.