

複数情報源に対する主成分分析

諏訪恭平／富岡亮太／矢入健久／○鹿島久嗣

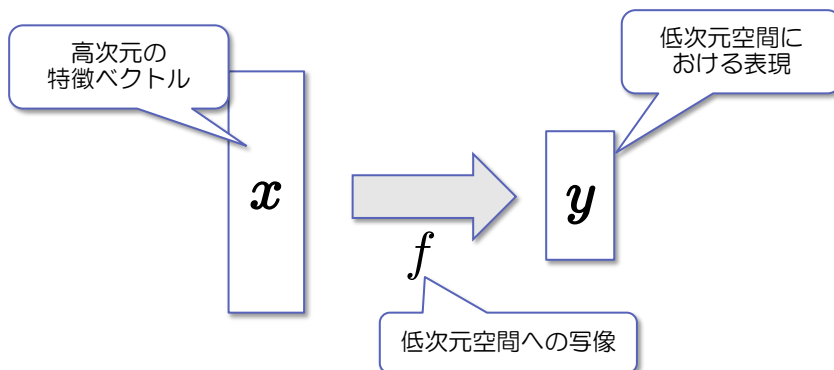


複数種類の情報源から得られたデータに対する 主成分分析法を提案します

- はじめての複数情報源に対する主成分分析法
 - 重要な情報源を自動的に選択する機構をもつ
- みどころ：
 - グループラッソ制約を用いることで情報選択機構をもつ
 - 主成分ごとに異なる情報源を選択する
 - グループラッソから複数カーネル学習の導出レシピ
- 観測データからの物体配置図復元タスクで有効性を検証
 - 相補的な情報を適切に組み合わせる
 - 不要な情報源を切り捨て、有用なもののみを残す
- いまのところ実験的にはいまひとつ...

次元削減は 高次元空間のデータを 低次元空間に落として
扱いやすくするための技術 です

- 多くの分析対象は高次元の実数値特徴ベクトル x として表現される
- もとのデータ x の性質を保ったまま、低次元の空間における y を得るための写像 f を求めるのが次元削減である
- 視覚化や前処理に利用される

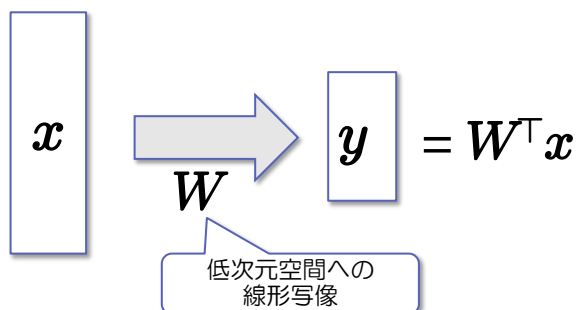


3

THE UNIVERSITY OF TOKYO

たとえば 主成分分析は 代表的な次元削減法のひとつ です

- 射影行列 W ($D \times d$ 行列) を用いて D 次元から d ($\ll D$)次元に射影
 - W : 正規直交基底 ($W^T W = I$)

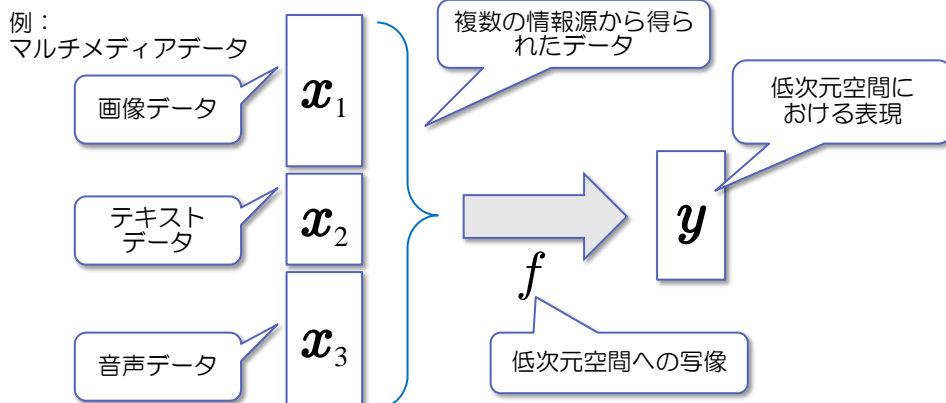


4

THE UNIVERSITY OF TOKYO

本研究では 複数の情報源からデータが得られるような場合の次元削減を考えます

- 多くの場合、一つの対象は複数の情報源からのデータをもつ
 - マルチメディアデータ：画像、テキスト、音声、...
 - 医療データ：検査項目、SNP、遺伝子発現、...

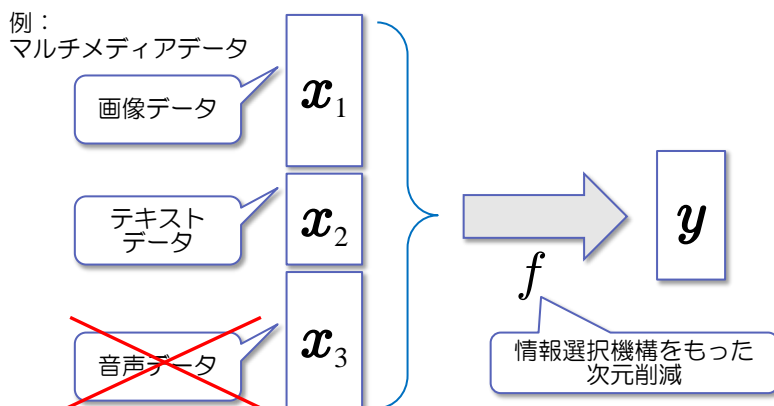


5

THE UNIVERSITY OF TOKYO

本研究では 複数の情報源の取捨選択/重みづけを自動的に行う次元削減法を考えます

- データには特徴的な構造をもったものとそうでないものがある
- これらを自動的に取捨選択または重みづけする次元削減法がほしい
- 具体的には主成分分析をベースに考える



6

THE UNIVERSITY OF TOKYO

主成分分析の復習

7

THE UNIVERSITY OF TOKYO

主成分分析の復習： 射影の分散を最大にするような写像を求めます

- 1次元に落とす場合を考える
 - X ：データのデザイン行列 (X^T の列が各特徴ベクトル)
 - w ：射影ベクトル
- 射影後の分散を最大にするような最適化問題を解くことで、射影ベクトル w が求まる

$$\text{maximize}_w w^T X^T X w \quad (\text{射影後の分散})$$

$$\text{s.t. } \|w\|_2^2 \leq 1 \quad (\text{射影ベクトルの大きさを制限})$$

8

THE UNIVERSITY OF TOKYO

主成分分析の復習：
最適な射影は固有値問題を解くことで得られます

- 最適解は固有値問題を解くことによって求まる

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \lambda \mathbf{w}$$

- 最大固有ベクトルを求める
- 2次元以上に落とす場合には：
 - 方法1：固有値の大きいほうから固有ベクトルを複数とる
 - 方法2： d 次元目までの寄与分を取り除いたデザイン行列
$$\mathbf{X}^{(d+1)} \equiv \mathbf{X}^{(d)} - \mathbf{y}^{(d)} \mathbf{w}^{(d)T}$$
に対し固有値問題の最大固有ベクトル $\mathbf{w}^{(d+1)}$ を求める操作を繰り返す

提案手法：複数情報源の主成分分析

提案手法：主成分分析の目的関数における2-ノルム制約をグループラッソの制約でおきかえます

- 主成分分析を情報源が M 個ある場合に拡張する
 - $\mathbf{X} \equiv (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_M)$: M 個の情報源データのデザイン行列
 - $\mathbf{w} \equiv (w_1, w_2, \dots, w_M)$: 射影ベクトル
- 最適化問題における2ノルム制約を、グループラッソで用いられる制約で置き換える

$$\text{maximize}_{\mathbf{w}} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} \quad (\text{射影後の分散})$$

$$\text{s.t. } \sum_m \|\mathbf{w}_m\|_2 \leq 1 \quad (\text{射影ベクトルの大きさを制限})$$

$$\text{※もともとの制約は } \sum_m \|\mathbf{w}_m\|_2^2 \leq 1$$

- グループラッソの制約は各情報源の射影ベクトル単位で0につぶす
 - 1-ノルムの特徴選択効果を特徴のグループに対して与える

最適化のテクニック：
グループラッソ制約を扱いやすい2ノルム制約に帰着します

- ラグランジアンはそのまま最適化しにくい 2-ノルムの1乗

$$L(\mathbf{w}, \lambda) \equiv \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \lambda (\sum_m \|\mathbf{w}_m\|_2 - 1)$$

- ラグランジアンを扱いやすい形に変形する

- 相加相乗平均から

$$\|\mathbf{w}_m\|_2 = \min_{\beta_m \geq 0} \frac{1}{2} \left(\frac{\|\mathbf{w}_m\|_2^2}{\beta_m} + \beta_m \right)$$

2-ノルムの2乗!
(扱いやすい)

- これを持ちいて

$$\bar{L}(\mathbf{w}, \beta, \lambda) \equiv \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \lambda \left(\sum_{m=1}^M \frac{1}{2} \left(\frac{\|\mathbf{w}_m\|_2^2}{\beta_m} + \beta_m \right) - 1 \right)$$

β が登場

- 「2ノルムの2乗」の項を「2ノルムの2乗」で置き換えた
- 但しパラメータ $\beta = (\beta_1, \beta_2, \dots, \beta_M)$ が増えていることに注意

新たに導入されたパラメータ β は各情報源の重要度として解釈できます

- $1/\beta_m$ は m 番目の情報源に対応する $\|w_m\|_2^2$ の係数になっている

$$\bar{L}(w, \beta, \lambda) \equiv wX^T Xw - \lambda \left(\sum_{m=1}^M \frac{1}{2} \left(\frac{\|w_m\|_2^2}{\beta_m} + \beta_m \right) - 1 \right)$$

β が大きい \Rightarrow
 w_m の正則化が弱い

- β_m は m 番目の情報源の重要度として解釈できる
 - β_m が大きいほど w_m の正則化が弱まる $\Rightarrow w_m$ を大きくできる
 - β_m が小さいほど w_m の正則化が強まる $\Rightarrow w_m$ が0になりやすい

逐次最適化アルゴリズム：
射影ベクトル w と重要度パラメータ β を交互に最適化します

- 射影ベクトル $w \equiv (w_1, w_2, \dots, w_M)$ についての最大化
 - $w \equiv (w_1, w_2, \dots, w_M)$ と $X \equiv (X_1, X_2, \dots, X_M)$ をスケーリングする

$$\Rightarrow \tilde{w}_m \equiv \frac{w_m}{\sqrt{\beta_m}} \quad \Rightarrow \tilde{X}_m(\beta_m) \equiv \sqrt{\beta_m} X_m$$

- 通常の主成分分析と同様の形のラグランジアンが得られる

$$\bar{L}(\tilde{w}, \lambda) = \tilde{w}^T \tilde{X}(\beta)^T \tilde{X}(\beta) \tilde{w} - \frac{\lambda}{2} \left(\|\tilde{w}\|_2^2 + \sum_{m=1}^M \beta_m - 2 \right)$$

- 固有値問題に帰着される

$$\tilde{X}(\beta)^T \tilde{X}(\beta) \tilde{w} = \frac{\lambda}{2} \tilde{w}$$

- 最大固有ベクトルを求める

- 重要度パラメータ $\beta = (\beta_1, \beta_2, \dots, \beta_M)$ についての最大化

- こちらは閉じた形で解が求まる $\beta_m^{new} = \sqrt{\beta_m^{old} \|\tilde{w}_m\|_2}$

主成分の逐次抽出：第 d 主成分までの寄与分をデータから取り除き、第 $d+1$ 主成分を求めます

- 繰り返しのよって、主成分が1軸求まる
 - 射影ベクトル w についての最大化（最大固有ベクトル）
 - 重要度パラメータ β についての最大化（閉じた形の解）
- さらにこれを繰り返すことで、1軸ずつ主成分を求める

第 d 主成分は $y^{(d)} \equiv \sum_{m=1}^M y_m^{(d)}$ $y_m^{(d)} \equiv X_m w_m^{(d)}$

- 第 d 主成分の寄与分をデータから取り除く

$$X_m^{(d+1)} \equiv X_m^{(d)} - y^{(d)} w_m^{(d)\top}$$

- 主成分ごとに異なった情報源重要度が得られる
 - 既存の情報源選択法と異なる点

マルチカーネル化：2ノルム制約が現れたことによりカーネル関数による非線形化が可能になります

- 目的関数の変形によって2ノルムの2乗が登場した

$$\tilde{L}(\tilde{w}, \lambda) = \tilde{w}^\top \tilde{X}(\beta)^\top \tilde{X}(\beta) \tilde{w} - \frac{\lambda}{2} \left(\|\tilde{w}\|_2^2 + \sum_{m=1}^M \beta_m - 2 \right)$$

- カーネル法による非線形化が可能
 - パラメータをおきかえる $\tilde{w}_m \equiv \tilde{X}_m(\beta)^\top \tilde{\alpha}$
 - 固有値問題は：

$$\tilde{K}(\beta) \tilde{\alpha} = \frac{\lambda}{2} \tilde{\alpha} \quad \text{ただし} \quad \tilde{K}(\beta) \equiv \sum_{m=1}^M \beta_m K_m$$

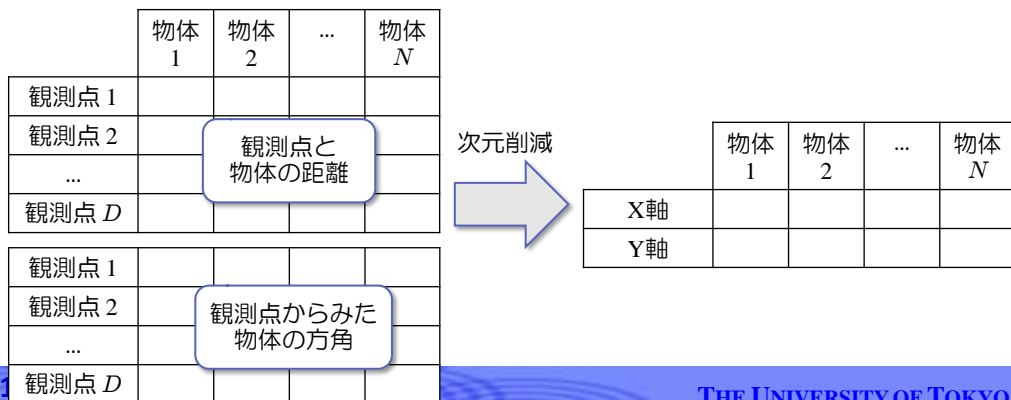
$$K_m \equiv X_m X_m^\top$$

情報源ごとの
カーネル関数
(マルチカーネル)

数値実験：物体配置図の復元タスク

物体配置図の復元タスクを題材に その復元精度で性能評価します

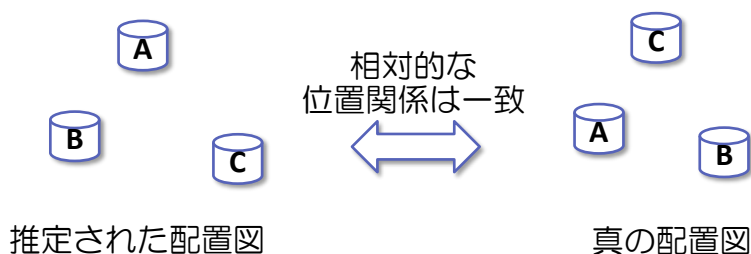
- 複数の物体の配置図を複数の観測点での観測データから復元する
 - 観測点と観測対象の座標はともに不明
- 例えば、各観測点からの各物体への距離や方角などが観測される
- 次元削減によって2次元に落とすことで配置図を復元する



評価方法：

配置図の復元精度は任意の3点の相対位置の一致度で測ります

- 教師なし次元削減では絶対的な座標は復元されない
 - わかるのは対的な位置関係のみ
- 復元精度は任意の3点の相対位置が一致する割合で評価する
(= 時計回りの順が一致する割合)

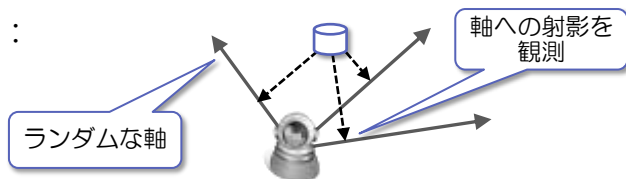


19

THE UNIVERSITY OF TOKYO

実験1：相補的な情報源の統合

- ランダムな3方向に物体を射影したときの距離を3つの情報源とする
 - 相補的に働く
- これらを組み合わせることで配置図の復元精度が上がるか？
- 比較手法：
 - 提案手法
 - 通常の主成分分析
 - 各情報源で主成分分析を行ったもののうち最も精度の高いもの
- ランダムにデータを生成：
 - 物体数：50
 - 観測点数：500

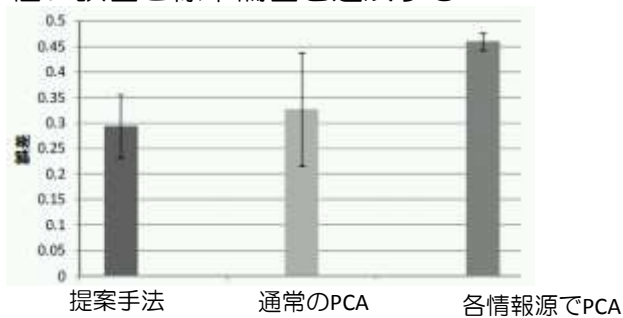


20

THE UNIVERSITY OF TOKYO

実験1：提案手法は相補的な情報をうまく組み合わせることが出来ます

- 比較手法
 - 提案手法
 - 通常の主成分分析
 - 各情報源で主成分分析を行ったもののうち最も精度の高いもの
- 提案手法が低い誤差と標準偏差を達成する



21

THE UNIVERSITY OF TOKYO

実験2：有用でない情報源の切り捨て

- ただ一つの有用な情報源に、役に立たない（ランダムな）情報源を加えて行った場合の復元精度への影響を調べる
- 有用な情報源として：
 - 観測点から一定距離内に物体がある（1）かない（0）か

	物体 1	物体 2	...	物体 N
観測点 1	1	0	...	1
観測点 2	1	1	...	0
...
観測点 D	0	0	...	1

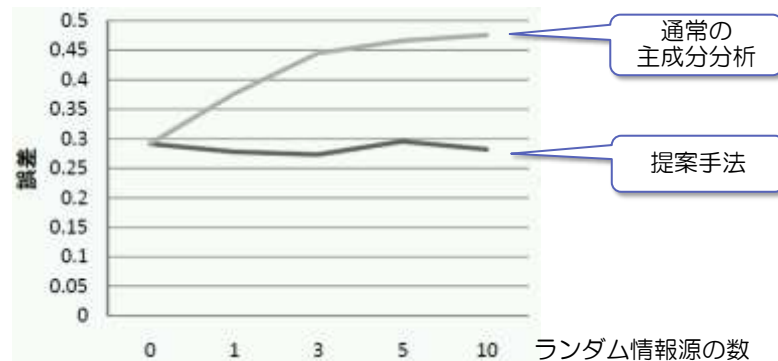
各観測点から一定距離内にその物体があるかどうか

22

THE UNIVERSITY OF TOKYO

実験2：提案手法はノイズに影響されず有用な情報源のみを選び出すことができます

- ランダム情報源の数を増やしていくと：
 - 通常の主成分分析では精度が悪化していく
 - 提案手法では悪化しない



23

THE UNIVERSITY OF TOKYO

まとめ: 複数情報源から得られたデータに対する主成分分析法を提案しました

- 複数情報源に対する主成分分析法を提案した
 - ポイント：
 - グループラッソ制約を主成分分析に用いることで情報源選択機構を実現する
 - 情報源の重要度を導入することで問題を解きやすくする
 - 主成分ごとに異なる情報源を選択する
 - 提案手法は：
 - 相補的に働く情報源を適切に組み合わせることができる
 - 不要な情報源を取り除くことができる
- ことを物体配置図の復元タスクで検証した

24

THE UNIVERSITY OF TOKYO