

局所線形アライメントによる欠損データ行列の共埋め込み法

Co-Embedding Method for Incomplete Data Matrix by Locally Linear Alignment

矢入健久*1

Takehisa Yairi

*1東京大学先端科学技術研究センター

RCAST, University of Tokyo

In this paper, we consider the “co-embedding” problem whose goal is to embed both the row and column vectors of a given observation matrix into low-dimensional latent spaces respectively. We propose an efficient co-embedding algorithm especially when a large portion of the matrix is structurally missing. A remarkable feature of the method is that the low-dimensional latent representations are efficiently obtained via matrix eigendecomposition, whereas conventional methods require iterative estimation of missing values and are subject to local optima. Besides, we extend the unsupervised co-embedding method to a semi-supervised version, in which the computation is reduced to solving a system of linear equations.

1. はじめに

次元削減は、分類、回帰、クラスタリングなどにも主要な機械学習技術であり、膨大かつ高次元なデータを可視化したり、潜在的次元を抽出する目的などで利用されている。

基本的な次元削減法は、多変量解析の古典的手法である主成分分析 (PCA) や多次元尺度構成法 (MDS) などに溯るが、近年では Isomap [Tenenbaum 00] や LLE [Roweis 00] に代表される非線形次元削減の研究が盛んである。また、関連分野として、特異値分解 (SVD) や非負行列分解 (NMF) など、高次元行列の低ランク近似に関する理論的研究、応用も数多く報告されている。

しかし、次元削減や行列分解を応用する際に問題となるのが、実世界のデータは、単に高次元かつ膨大だけでなく、様々な理由により、欠損を伴うという点である。これに対して、上述の次元削減・行列分解の手法は完全なデータを想定しているため、そのままでは適用することができない。

データの欠損に対する最も単純なアプローチは、欠損値を定数で埋めることである。例えば、欠損が「ほぼ 0」を意味するとか、センサーの計測範囲外を表すようなデータを扱う場合、このアプローチは合理的である。しかし、全てのデータがそのような性質を持つわけではなく、どのような定数を用いるかが自明でない場合も多い。

一方、より洗練された方法は、欠損値の推定と次元削減とを交互かつ反復的に行うことである。そのようなアルゴリズムとしては、交互最小二乗法 (ALS) [Shum 95] や Wiberg アルゴリズム [Okatani 07] などが知られており、画像列から対象物体の 3 次元モデルを推定する Structure from Motion (SFM) などに用いられている。また、機械学習分野でも、確率的 PCA (PPCA) を EM アルゴリズムを用いて学習することにより、欠損値を推定しながら次元削減を行うことができるが示されている [Roweis 98]。しかし、これらの反復推定アプローチは、欠損が大規模である場合には、収束性も解の質も劣化する。さらに、欠損に構造的なパターンがある場合、すなわち、個々の変数が欠損する確率が、求めるべき潜在変数ベクトルの値に依存する場合、その情報を利用するメカニズムを持たない。

これに対して、提案手法は、欠損値を定数で埋めることも、

反復的に推定することもしない。代わりに、欠損していない値だけを利用する。提案手法は、非欠損要素に関して、求めるべき列潜在ベクトルとの間に大雑把に線形性が成り立つという仮定に基づいて、欠損パターンの構造を積極的かつ有効に利用する。その結果、大規模かつ構造的に欠損している行列データの行ベクトル集合および、列ベクトル集合、それぞれの低次元潜在表現を求めることができる。

2. 局所アライメントによる共埋め込み法

2.1 構造的欠損を伴う行列データの共埋め込み問題

$M \times N$ の観測データ行列 $Y = [y_{i,j}]_{i=1,\dots,M,j=1,\dots,N}$ を考える。ただし、各 (i,j) 成分 $y_{i,j}$ はスカラー値に限定されず、一般に D 次元ベクトルであるとする。*1

Y の (i,j) 成分 $y_{i,j}$ が存在しているか欠損しているかを表すために、インジケータ変数 $q_{i,j}$ を導入する。すなわち、もし $y_{i,j}$ が存在していれば $q_{i,j} = 1$ 、欠損していれば $q_{i,j} = 0$ とする。

今、共埋め込み (Co-embedding) は、以下の 2 つの問題を同時に解くこととして定義される。

1. Y を構成する M 個の ND 次元の行ベクトル集合を次元削減して、 n 次元ベクトルの集合 $X = [x_1, \dots, x_M]^T$ を得ること。
2. Y を構成する N 個の MD 次元の列ベクトル集合を次元削減して、 m 次元ベクトルの集合 $Z = [z_1, \dots, z_N]^T$ を得ること。

ここで注意すべき点は、行列の低ランク近似とは異なり、 X と Z の積によって、(欠損の無い) Y を復元することが目的なわけではないことである。

この行列データの共埋め込み問題は、以下のような解釈も可能である。ある観測 $y_{i,j}$ が、2 つの潜在変数ベクトル $x_i \in \mathcal{R}^n$ と $z_j \in \mathcal{R}^m$ に関する「未知の」2 項関数 g によって生成されると考える。すなわち、

$$y_{i,j} = g(x_i, z_j) + e_{i,j} \quad (1)$$

*1 したがって、 Y は厳密にはテンソルであるが、本論文では行列の拡張として考える方が理解しやすいので、敢えて行列として扱う。

ここで、 $e_{i,j}$ は観測ノイズを表す。このとき、共埋め込みの目的は、一部の (i, j) の組み合わせについて与えられた $y_{i,j}$ の集合から、すべての $\{x_i\}$ ($i = 1, \dots, M$) および $\{z_j\}$ ($j = 1, \dots, N$) を求めることである。

以上が一般的な欠損を伴う行列の共埋め込み問題であるが、本論文では対象を限定し、欠損が構造的パターンを持つ場合を考える。すなわち、ある観測データ $y_{i,j}$ の存在性が、 z_j に関して局所性を持つという仮定を置く。言い換えれば、「ある j, j' に対して、 $q_{i,j} = q_{i,j'} = 1$ となる i が存在するとき、 z_j と $z_{j'}$ は近傍にある」という仮定である。

一見、この仮定はかなり限定的なものであるように思えるが、後の実験の章で述べるように、このような性質が近似的に成り立つ実問題は多く存在する。

2.2 局所線形アライメントの基本アイデア

上の仮定が成り立つとすると、 $q_{i,j} = 1$ のとき、すなわち、 $y_{i,j}$ の値が存在するとき、その近傍で近似的に局所線形性が成り立つと考えられる。すなわち、

$$y_{i,j} = g(x_i, z_j) + e_{i,j} \approx G(x_i)[z_j^\top, 1]^\top = G(x_i)\tilde{z}_j \quad (2)$$

ここで、 $G(x_i)$ は x_i によって決まる射影行列であり、 \tilde{z}_j は z_j の同次座標を表すとす。

この近似は次のように解釈することができる。今、 x_i は時刻 i における観測者の潜在状態（例えば位置と姿勢）を表し、 z_j は j 番目の物体の潜在状態（例えば位置）を表すと考える。すると、式 2 は、「観測者の状態 x_i によって定まる観測部分空間に、その状態で観測することができる物体の状態 z_j を線形射影したものが、観測データ $y_{i,j}$ である」ということを表している。言い換えれば、「各時刻における観測データは、その場で見える世界の一部分を低次元の知覚空間に線形射影したものとみなせる」ということである。

共埋め込みの第一の目的は、観測データの断片群をうまく整列（アライメント）して、物体集合の潜在状態 $\{z_j\}$ 、すなわち、世界を復元することである。また、各観測断片に対するアライメント操作は、その断片を生成した観測者の状態を反映していると考えられるので、それらのアライメントのデータを次元削減することによって、観測者の状態 $\{x_i\}$ も復元される。

2.3 教師なし局所線形アライメント共埋め込み

まず、列潜在ベクトル集合 $\{z_j\}$ の復元を考える。前節で述べた仮定は、もし、 $q_{i,j} = 1$ であれば、 z_j が $y_{i,j}$ に対して近似的に線形（より厳密にはアフィン）であることを意味する。そこでこの仮定を利用し、ある z_j について、 $y_{i,j}$ をアフィン写像したのによって逆に近似することを考える。すなわち、

$$\tilde{z}_{i,j} \equiv T_i[y_{i,j}^\top, 1]^\top = T_i\tilde{y}_{i,j} \quad (3)$$

ここで、変換行列 T_i はアライメント操作を表す行列であり、 $q_{i,j} = 1$ が成り立つ全ての j について共通である。ただし、 $\tilde{y}_{i,j}$ は $y_{i,j}$ の同次座標表現とする。このとき、最終的な z_j の推定値は、 $q_{i,j} = 1$ であるような全ての i について $\tilde{z}_{i,j}$ の平均によって得られると考える。

$$\hat{z}_j = \frac{\sum_{i=1}^M q_{i,j} \tilde{z}_{i,j}}{\sum_{i=1}^M q_{i,j}} = \sum_{i=1}^M \tilde{q}_{i,j} \tilde{z}_{i,j} \quad (4)$$

ただし、 $\tilde{q}_{i,j} = q_{i,j} / \sum_{i=1}^M q_{i,j}$ は規格化されたインジケータ変数である。

さて、ここでの最大の問題は、いかにして変換行列の集合 $\{T_i\}$ ($i = 1, \dots, M$) を求めるかということである。一つの合理的な方法は、 z_j の推定値である $\{\hat{z}_{i,j}\}$ ($i = 1, \dots, M$) ができるだけ互いに一致するように、各 T_i を選ぶことである。具体的には、以下の式で定義されるコスト関数 Φ_{aln} が最小となる $\{T_i\}$ ($i = 1, \dots, M$) を求める。

$$\Phi_{aln} = \frac{1}{2} \sum_{j=1}^N \sum_{i \neq i'} \tilde{q}_{i,j} \tilde{q}_{i',j} \|\hat{z}_{i,j} - \hat{z}_{i',j}\|^2 \quad (5)$$

スペースの都合上、詳細な導出は割愛するが、式 5 は、次式のように行列の二次形式のトレースとして書くことができる。

$$\Phi_{aln}(T) = \text{Tr}(T^\top (D - V^\top V) T) \quad (6)$$

ただし、この式では補助的な行列およびベクトルを以下のように定義している。 $v_j = [\tilde{q}_{1,j} \tilde{y}_{1,j}^\top, \dots, \tilde{q}_{M,j} \tilde{y}_{M,j}^\top]$ 、 $V = [v_1^\top, \dots, v_N^\top]^\top$ 、 $D_i = \sum_{j=1}^N \tilde{q}_{i,j} \tilde{y}_{i,j} \tilde{y}_{i,j}^\top$ 、 $D = \text{diag}(D_1, \dots, D_M)$ 、 $T = [T_1, \dots, T_M]^\top$ 。また、求めるべき $Z = [z_1, \dots, z_N]^\top$ は、 $Z = VT$ と書ける。

Φ_{aln} の最小化は、 $T = 0$ という自明な解が存在するので、 $Z^\top Z = T^\top (V^\top V) T = I$ という制約を付けて最小化する。この制約付き最小化問題の解は、 $T_{opt} = [u_2, \dots, u_{m+1}]$ として得られる。ただし、 u_2, \dots, u_{m+1} は、一般化固有値問題

$$(D - V^\top V)u = \lambda(V^\top V)u \quad (7)$$

の 2 番目に小さい固有値から $m+1$ 番目に小さい固有値に対応する固有ベクトルを並べたものである。また、 \hat{Z} は、 $\hat{Z} = VT_{opt}$ として得られる。

次に、行潜在ベクトル集合 $X = [x_1, \dots, x_M]^\top$ の復元を考える。上で求めた各アライメント行列 T_i は、対応する行潜在ベクトル x_i に依存していると考えられる。したがって、 $\{x_i\}$ は、 $\{\text{vec}(T_i)\}$ を n 次元に削減することによって得られると期待される。ただし、 $\text{vec}(T_i)$ は行列 T_i に含まれる全列ベクトルを縦に並べて得られるベクトルを表す。もとの観測行列 Y は多くの欠損値を含んでいたのに対して、 T には欠損が無いことに注意されたい。具体的な次元削減法としては、非線形のものも含めて様々考えられるが、今回は最も単純な特異値分解 (SVD) を用いた。

なお、上のコスト関数および列潜在ベクトル z_j の求め方は、基本的に Verbeek らによる非線形 CCA [Verbeek 04] の解法に基づいている。ただし、彼らの手法は、非線形次元削減の一手法として提案されたものであり、欠損データを扱うものではない。また、行と列の同時次元削減の考え方はないという点で本提案手法とは異なる。

2.4 正則化と半教師あり共埋め込み

前節で示した共埋め込みアルゴリズムは、求めるべき潜在変数ベクトル集合 $\{x_i\}$ 、 $\{z_j\}$ に関して事前知識を想定しないという点で、「教師なし」の手法と言える。一方、実際の応用では、分野固有の様々な事前知識を用いることによって、共埋め込みの性能を向上させることが可能である。

例えば、行潜在ベクトル x_i が、ゆっくりと移動する観測者の時刻 i における位置・姿勢を表す場合、次時刻の行潜在ベクトル x_{i+1} と大きくは異ならないと期待される。この正則化は次式で表されるコスト関数によって表現できる。

$$\Phi_{smo} = \text{Tr}(T^\top (S^\top S) T) \quad (8)$$

ただし、行列 S は、 T の連続する要素間の差を計算する演算子を表す。このコスト関数を重み付きで Φ_{aln} に加えたもの、すなわち、 $\Phi_{aln} + \alpha_{smo} \cdot \Phi_{smo}$ の最小化は同様に固有値問題に帰着される。

さらに、アプリケーションによっては半教師ありの問題、すなわち、一部の行・列潜在ベクトルについて、正解（ラベル）が事前に与えられる状況が考えられる。例えば、前述の移動観測の場合、一部の観測位置やランドマークの位置が GPS 計測などによって得られることがある。提案手法はこのような半教師ありの問題設定にも容易に拡張することができる。

まず、列潜在ベクトル z_j のラベルデータを z_j^* と表し、その値が与えられるか否かを示す二値変数 δ_j を定義する。

$$z_j^* = z_j \text{ (if } \delta_j = 1\text{)}, \quad \mathbf{0} \text{ (if } \delta_j = 0\text{)} \quad (9)$$

このとき、求める列潜在ベクトルがラベルデータとどれだけ一致するかを表すコスト関数は以下のように書ける。

$$\Phi_{zlb} \equiv \sum_{j=1}^N \delta_j \|z_j - z_j^*\|^2 \quad (10)$$

ここで、 $Z^* = [z_1^*, \dots, z_N^*]^\top$ および $J_z = \text{diag}(\delta_1, \dots, \delta_N)$ とすると、式 10 は、次式のように書ける。

$$\Phi_{zlb} = \text{Tr}((VT - Z^*)^\top J_z (VT - Z^*)) \quad (11)$$

これを重み付きで加えた全体のコスト関数 $\Phi_{sem}(T) = \Phi_{aln} + \alpha_{smo} \cdot \Phi_{smo} + \alpha_{zlb} \cdot \Phi_{zlb}$ の最小化は、線形の連立方程式を解くことに帰着される。

$$T_{opt} = (D + V^\top (\alpha_{zlb} J_z - I) V + \alpha_{smo} S^\top S)^{-1} (\alpha_{zlb} V^\top J_z Z^*) \quad (12)$$

次に行潜在ベクトル $\{x_i^*\}$ の一部についてラベルデータが与えられた場合を考える。これは、ベクトル化された T_i 、すなわち、 $\{\text{vec}(\hat{T}_i)\}$ を入力ベクトルとし、一部のターゲット変数（ベクトル）が $\{x_i^*\}$ として与えられるような、一般的な半教師あり回帰問題である。この問題に対しては、様々なアルゴリズムの適用が考えられるが、今回は、Tikhov 正則化項を加えた半教師あり線形回帰を用いた。

3. 実験

3.1 (実験 1) SFM 問題

最初の実験では、従来手法との比較を目的として、以下のような小規模な SFM (Structure from Motion) 問題、すなわち、画像列から被写体の 3 次元形状とカメラの姿勢を同時復元する問題を取り上げる。ここでは、正 12 面体の物体をランダムな方向から見て各頂点を識別し、それらのカメラ画像上の 2 次元座標が観測データ $\{y_{i,1}, \dots, y_{i,N}\}$ として得られるとする。ここで、正 12 面体は 20 個の頂点を持つので、 $N = 20$ となる。そして、 $M = 100$ 回の観測によって、観測行列 Y が得られたとする。このとき、3 次元モデル、すなわち 20 個の頂点の 3 次元座標 $\{z_j\}$ を復元することが目的である。

まず、常に全頂点が見えると仮定した場合、この問題は特異値分解 (SVD) によって解くことができる (図 1(a))。一方、提案手法 (教師なしバージョン) を適用した結果が図 1(b) であり、同様に正 12 面体の 3 次元形状が完全に復元された。

次に、遮蔽される頂点は観測されないという、現実的な条件で実験を行った。このとき、 Y の要素の約 30% が欠損した。こ

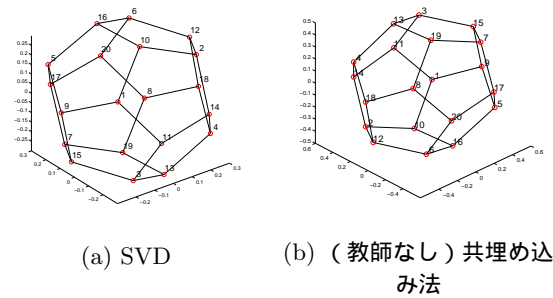


図 1: 欠損の無い観測データから再構成された 3 次元形状

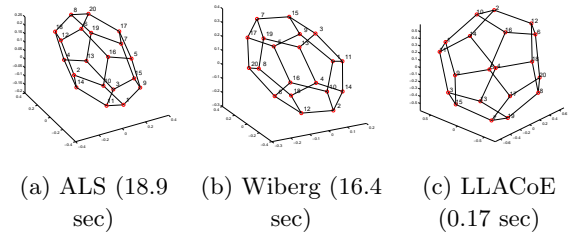


図 2: 欠損のある観測データから復元された正 12 面体の 3D モデルと計算時間. 全てのアルゴリズムは Matlab で実装され、Dell Precision T1500 で実行された

の場合、欠損値を定数値で代替するのは不適当であり、通常の特異値分解は利用できない。そこで、交互最小二乗法 (ALS)、および、Wiberg アルゴリズムを提案手法と比較した。図 2 (a)-(c) は、得られた 3 次元形状と、計算に要した時間を示したものである。いずれの三手法もほぼ正確に復元しているが、提案手法は従来手法と比べて計算時間が極めて少なくなっている。

3.2 (実験 2) SLAM 問題

次に、移動ロボット研究における SLAM (Simultaneous Localization and Mapping) 問題と同様のタスクに適用した。このタスクでは、仮想的なキャンパスに、564 個のアクセスポイント (AP) が散在し、観測者が通路上を移動しながら、観測地点で認識することのできた AP の大雑把な相対方向と相対距離を計測できると仮定する。図 3 は、このシミュレーションで想定したキャンパスの地図と、ランドマークの配置を示したものである。また、図 4 は観測者の実際の移動軌跡と 310 点

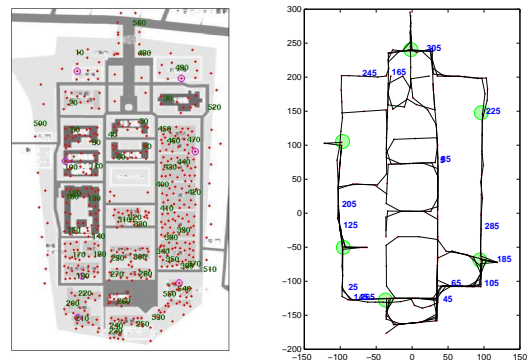


図 3: シミュレーション環境 図 4: 観測者の移動軌跡

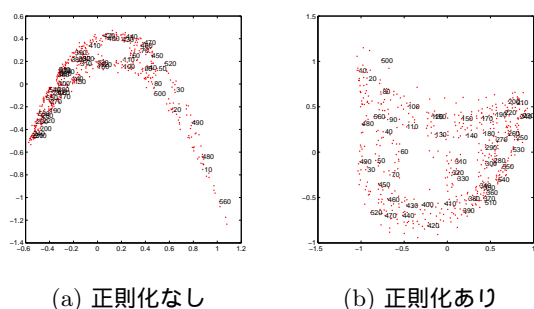


図 5: 教師なし問題設定下での推定された地図

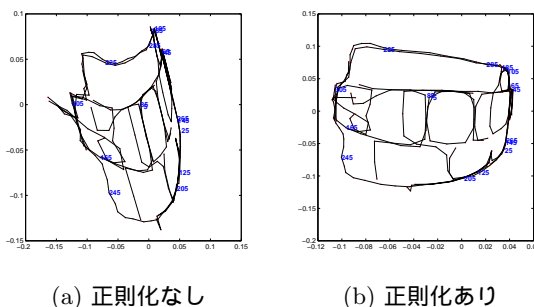


図 6: 教師なし問題設定下で推定された軌跡

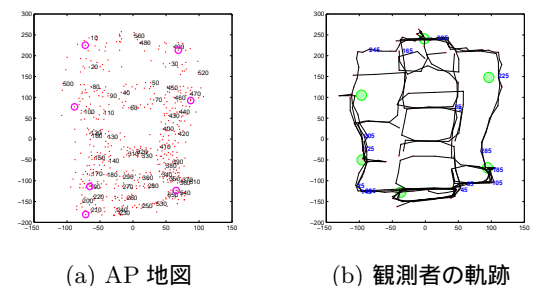


図 7: 教師あり共埋め込み法により推定された地図と軌跡

の観測点を図示したものである。

このタスクでは、行潜在ベクトル x_i ($i = 1, \dots, 310$) は移動観測者の状態 (すなわち、位置と向き) であり、列潜在ベクトル z_j ($j = 1, \dots, 564$) は AP の位置である。また、観測データ $y_{i,j}$ は、時刻 i における観測者から j 番目の AP への方角 (bearing) と距離 (range) から計算された相対位置とする。ただし、観測には大きなノイズが含まれる。

今回は、ある観測データ $y_{i,j}$ が得られる確率が次式で得られるとしてシミュレーションを行った。

$$Pr(q_{i,j} = 1) = \frac{1}{1 + \exp(0.15 \cdot (d_{i,j} - 50))} \quad (13)$$

ここで、 $d_{i,j}$ は時刻 i 番目の観測地点と j 番目の AP との距離である。シミュレーションの結果、 Y 中の欠損値の割合は約 97% に上った。

3.2.1 教師なし条件での地図・自己位置推定

まず、 X 、 Z のいずれについてもラベルデータが得られないという教師無しの条件、かつ、移動軌跡の連続性を表す正則化項 Φ_{smo} を用いないで提案手法を適用した。図 5(a) に示したように、AP の地図 (Z) は大きく歪んでしまっているが、大まかな相対的位置関係はほぼ復元されている。また、観測者の移動軌跡 (X) は、図 6(a) に示すように良く復元されている。

一方、正則化項 Φ_{smo} を加えた場合の結果が Fig.5(b) および Fig.6(b) である。ただし、正則化項に対する重みパラメータは $\alpha_{smo} = 0.2$ とした。図から明らかのように、精度が大きく改善されていることが分かる。

3.2.2 半教師あり条件での地図・自己位置推定

次に、半教師ありの問題設定において提案手法を適用した。すなわち、図 3 で丸印で強調されている 7 個の AP の正確な位置を $\{z_j\}$ に関するラベルデータとして与え、また、図 4 で緑色の丸で示された 6 個の領域に含まれる観測地点における正確な位置と姿勢が $\{x_i\}$ のラベルデータとして与えられるとした。

図 7 (a) および (b) は、提案手法によって推定された AP の地図、および、観測者の移動軌跡を示したものである。ラベル情報の効果により、回転および鏡像に関する任意性が排除され、真の地図、移動軌跡との整合性が向上している。

4. おわりに

本稿では、要素の大部分が構造的に欠損しているような行列データの行と列をそれぞれ低次元の潜在空間に埋め込む共埋め込み法を提案した。提案手法は、膨大な欠損値を反復的に推定する必要が無く、データ点数に比例するサイズの対称行列の固有値分解に帰着されるという長所を持つ。また、半教師ありの問題設定下では線形方程式に帰着され、さらに計算が容易になることを示した。

実験では、2 種類のタスク、SFM と SLAM を模擬したシミュレーションデータに提案手法を適用し、提案手法の有効性、他手法に対する優位性を示した。今後は、様々な分野の多くの問題への適用を検討する予定である。

参考文献

- [Okatani 07] Okatani, T. and Deguchi, K.: On the Wiberg Algorithm for Matrix Factorization in the Presence of Missing Components, *International Journal of Computer Vision*, Vol. 72, No. 3, pp. 329–337 (2007)
- [Roweis 98] Roweis, S.: EM Algorithms for PCA and SPCA, in *Advances in Neural Information Processing Systems*, pp. 626–632 (1998)
- [Roweis 00] Roweis, S. and Saul, L.: Nonlinear dimensionality reduction by locally linear embedding, *Science*, Vol. 290, pp. 2323–2326 (2000)
- [Shum 95] Shum, H. Y., Ikeuchi, K., and Reddy, R.: Principal Component Analysis with Missing Data and Its Application to Polyhedral Object Modeling, *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 17, No. 9, pp. 854–867 (1995)
- [Tenenbaum 00] Tenenbaum, J. B., Silva, V. D., and Langford, J. C.: A global geometric framework for nonlinear dimensionality reduction, *Science*, Vol. 290, pp. 2319–2323 (2000)
- [Verbeek 04] Verbeek, J., S.T.Roweis, , and N.Vlassis, : Non-linear cca and pca by alignment of local models, in *Procs. of NIPS* (2004)