

階層的GTMに基づく教師付き学習モデル

02103750 筑波大学 後藤正輝 GOTO Masateru

01105930 筑波大学 *香田正人 KODA Masato

1 はじめに

階層的GTM[2]はGTMをノードを持つツリーとして構成される, データ可視化のための教師なし学習モデルである. またGTM(Generative Topographic Mapping,[1])は, 高次元のデータの確率分布を, いくつかの潜在変数に置き換えて表現する確率モデルである.

本研究では, 階層的GTMが定める潜在変数の確率分布に基づいた教師付き学習モデルを提案する.

2 Generative Topographic Mapping

D 次元データ x の確率分布 $p(x)$ を L 個の潜在変数 u による表現に置き換えるため, 潜在変数空間上の点 u をデータ空間上の対応する点 y に写す関数 $y(u; W)$ を定める. このとき $y(u; W)$ は潜在変数空間をデータ空間に埋め込まれた L 次元多様体に写す (以下 $L=2$ とする).

潜在変数空間上に確率分布 $p(u)$ を定義すると, データ空間上に対応する確率分布 $p(x|W)$ が生じる. データが多様体上に厳密に分布する状況は非現実的であるので, ノイズを含めた x の分布を定義する.

$$p(x|u, W, \beta) = \left(\frac{\beta}{2\pi}\right)^{D/2} \exp\left\{-\frac{\beta}{2}\|y(u; W) - x\|^2\right\}$$

データの分布 $p(x|W, \beta)$ は, 潜在変数 u についての積分として次式で表される.

$$p(x|W, \beta) = \int p(x|u, W, \beta)p(u)du$$

u についての積分は一般には困難であるため, 潜在変数空間上に規則的に格子点 $u_i, i = 1, \dots, K$ を配置し, u_i を中心とするデルタ関数の和として $p(u)$ を定義する.

$$p(u) = \frac{1}{K} \sum_{i=1}^K \delta(u - u_i)$$

このとき $p(x|W, \beta)$ は次式で得られる.

$$p(x|W, \beta) = \frac{1}{K} \sum_{i=1}^K p(x|u_i, W, \beta)$$

データ $\{x_1, \dots, x_N\}$ が与えられると, EM アルゴリズムによって W, β の最尤推定値が得られる [1].

$y(u; W)$ として, 一般化線形モデルを用いる.

$$y(u; W) = W\phi(u)$$

ϕ の要素は, M 個の基底関数 $\phi_j(u), j = 1, \dots, M$ であり, 潜在変数空間上に格子状に配置されたガウス分布を用いる.

x_n が与えられると潜在変数の事後分布が生成される.

$$R_{in} = p(u_i|x_n, W, \beta) = \frac{p(x_n|u_i, W, \beta)}{\sum_{i'=1}^K p(x_n|u_{i'}, W, \beta)}$$

全データをまとめて可視化するために, 期待値を取ることににより $p(u_i|x_n, W, \beta)$ を要約する.

$$E[u|x_n, W, \beta] = \int p(u|x_n, W, \beta)u du = \sum_{i=1}^K R_{in}u_i$$

3 階層的GTM

階層的GTM T は, 各ノードがGTM M で構成されるツリーである. T によって定義される x の確率分布は次式で表される.

$$p(x|T) = \sum_{M \in \text{Leaves}(T)} \pi(M)p(x|M)$$

ここで $p(x|M)$ は $p(x|W_M, \beta_M)$ の省略表記である. また $\pi(M)$ はモデル M の混合係数であり, 次式で再帰的に定められる.

$$\begin{cases} \pi(\text{Root}) = 1 \\ \pi(M) = \prod_{i=2}^{\text{Level}(M)} \pi(\text{Path}(M)_i | \text{Path}(M)_{i-1}) \end{cases}$$

ただし $\text{Path}(M)_i$ はルートノードから M に至るパスを構成する i 番目のノードを示す.

x が与えられると, 潜在変数の事後分布が生成される.

$$p(u_i^M|x, T) = \sum_{M \in \text{Leaves}(T)} p(M|x)p(u_i^M|x, M)$$

4 階層的GTMに基づく教師付き学習

訓練データ $D = \{(x_1, y_1), \dots, (x_N, y_N)\}$ を学習するため, はじめに T のリーフノードに対応するGTM M の, 格子点 u_i^M についての y の期待値を次式で推定する.

$$E[y|u_i^M, M, D] = \frac{\sum_{n=1}^N y_n \cdot p(u_i^M|x_n, M)}{\sum_{n=1}^N p(u_i^M|x_n, M)}$$

データ x が入力されると, 潜在変数の分布 $p(u_i^M|x, T)$ が生じる (図1). このとき対応する出力 y の期待値を次式で与える.

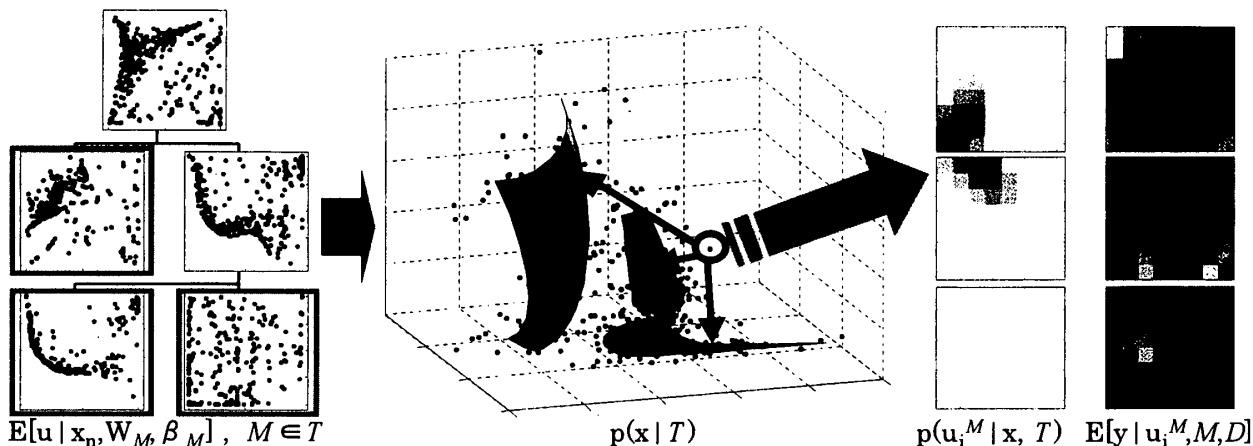


図 1: 階層的 GTM T は、リーフノードに対応する GTM M を混合係数 $\pi(M)$ で結合することにより、確率分布 $P(x|T)$ を構成する。データ x が与えられると、潜在変数の事後分布 $p(u_i^M|x, T)$ が生成される。また提案する教師付き学習モデルによって得られた $E[y|u_i^M, M, D]$ を全ての u_i^M についてプロットすることによって、 y の期待値のマップが生成される。

$$E[y|x, T, D] = \sum_{M \in \text{Leaves}(T)} E[y|x, M, D] p(M|x)$$

ただし $E[y|x, M, D]$ は次式で定められる。

$$E[y|x, M, D] = \sum_{i=1}^K E[y|u_i^M, M, D] p(u_i^M|x, M)$$

階層的 GTM T のリーフノードに対応する GTM M それぞれをクラスターとみなし、ファジィ・クラスタリングを規定する。データ (x, y) のクラスター M への所属度を表すメンバーシップ関数 μ_M を次式で定義する。

$$\mu_M(x, y) = P(M|x), \quad M \in \text{Leaves}(T)$$

また、クラスター M に属するデータの集合 C_M を次式で定める。

$$C_M = \{(x, y) \mid M = \arg \max_N \mu_N(x)\}$$

5 モデル構築の自動化

階層的 GTM T の構築過程で、GTM M に属するデータ C_M の分割 C_{M_L}, C_{M_R} を定める。

$$C_{M_L}(i, j) = \{(x, y) \in C_M \mid d(x, W_M \phi_M(u_i^M)) \leq d(x, W_M \phi_M(u_j^M))\}$$

$$C_{M_R}(i, j) = C_M - C_{M_L}(i, j)$$

ここで $d(\cdot, \cdot)$ は 2 点間のユークリッド距離を示す。 C_{M_L}, C_{M_R} それぞれに主成分分析を適用し、その結果を近似するように W, β の EM アルゴリズムの初期値を定めることによって、 M を 2 つの子ノードに分割する。

ノード分割の効果を評価するために、ノードに属するデータの逸脱度 $D(A)$ を定義する。教師信号が連続値の場合、平均二乗誤差を用いる。

$$D(A) = \frac{1}{|A|} \sum_{(x, y) \in A} (y - \mu_A)^2$$

また教師信号がバイナリの場合、エントロピーを用いる。

$$D(A) = -\mu_A \log \mu_A - (1 - \mu_A) \log(1 - \mu_A)$$

ここで μ_A は A に属するデータの平均であり、 $|A|$ は A に属するデータの件数である。分割による $D(A)$ の減少量を分割の効果 ΔD として定め、この値を最大にする i, j に基づいて子モデルの初期化を行う。

$$\Delta D(C_{M_L}, C_{M_R}) = D(C_M) - \left\{ \frac{|C_{M_L}| D(C_{M_L}) + |C_{M_R}| D(C_{M_R})}{|C_M|} \right\}$$

ノード分割の停止基準として、 $|C_M|$ が $|D|$ の一定以下になったら分割を停止する minSize と、 $|C_M| \cdot D(C_M)$ が $|D| \cdot D(D)$ の一定以下になったら分割を停止する minDev の 2 種類を定める。

6 数値実験

数値実験結果について当日報告する。

参考文献

1. C.M. Bishop, and M. Svensén and C.K.I. Williams. GTM: The Generative Topographic Mapping. *Neural Computation*, Vol. 10, No. 1, pp. 215-234, 1998.
2. P. Tino and I. Nabney. Hierarchical GTM: constructing localized non-linear projection manifolds in a principled way. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 24, No. 5, pp. 639-656, 2002.