

# モード探索型クラスタリング

江口 真透 数理・推論研究系 教授

**【アブストラクト】**

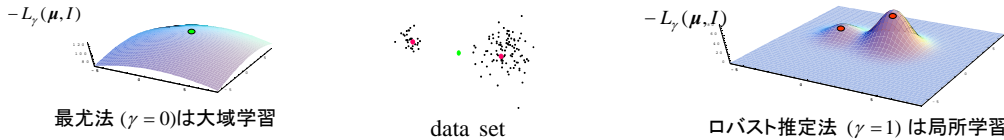
非階層クラスタリングの方法を提案する。従来のクラスター法は、例えば  $k$ -平均法、ファジィC-平均法、ガウス混合モデル法などはクラスター数を予め決めないと実行できない弱点がある。この発表では、**クラスター数を逐次的に決める方法**を提案する。タリス・エントロピーに関連したパワー・ロス関数を計算して、その局所最小解の個数を  $k$  とする。このように  $k$  を自動的に決め、 $k$  平均(局所最小解)からの距離によって  $k$  個のクラスターを定めることができる。パワー・エントロピーがデータセットを異なる定常状態をもつ  $k$  個のサブセットに分解していると解釈できる。

**【動機】**

$p$ 次元のデータ  $x_1, \dots, x_n$  に対して正規モデル  $N(\mu, V)$  に対して次のロス関数を考察していた。

$$L_\gamma(\mu, V) = -\det(V)^{-\frac{\gamma}{1+\gamma}} \sum_{i=1}^n [\exp\{-\frac{\gamma}{2}(x_i - \mu)^T V^{-1}(x_i - \mu)\} - 1] / \gamma \quad (1)$$

以下のような2次元データに対して平均ベクトル  $\mu$  の推定をロス関数 (1) の最小化で行うと奇妙なことが...



**【方法】**

予め分散行列  $V_0$  を決めて、 $L_\gamma(\mu, V_0)$  の局所最小解  $\mu_1, \dots, \mu_k$  を  $k$ -平均と定義して、 $k$ -クラスターを次のように定める。

$$C_j = \{ \mu : (\mu - \mu_j)^T V_0^{-1}(\mu - \mu_j) \leq (\mu - \mu_\ell)^T V_0^{-1}(\mu - \mu_\ell) \ (\forall \ell \neq j) \} \quad (j=1, \dots, k) \quad (2)$$

**【アルゴリズム】**

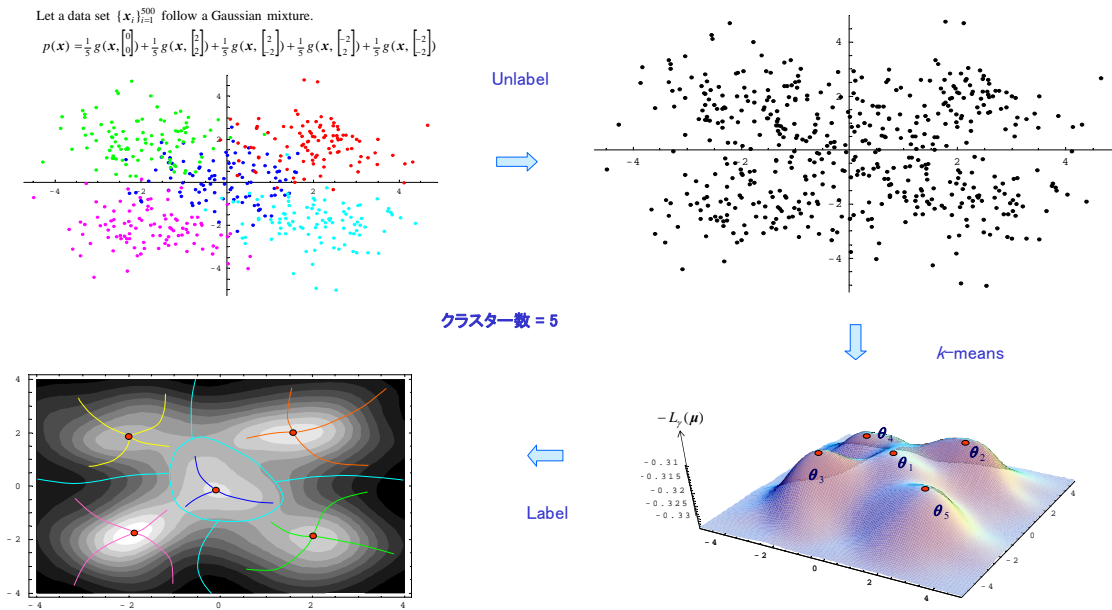
(A1) データ  $\{x_1, \dots, x_n\}$  の中からロス関数  $L_\gamma$  を最小にするものを初期値として、反復(\*)の収束値  $\mu_1$  を求める。

(A2)  $t=2, \dots$  に対して、データ  $\{x_1, \dots, x_n\}$  の中から  $\{\mu_1, \dots, \mu_{t-1}\}$  への距離:  $\delta(x) = \min_{\mu \in \{\mu_1, \dots, \mu_{t-1}\}} (x - \mu)^T V_0^{-1}(x - \mu)$  の  $\alpha$  番目に大きいものを初期値として反復(\*)の収束値  $\mu_t$  を求める。

$$(*) \quad \mu \leftarrow \frac{\sum_{i=1}^n w_i(\mu) x_i}{\sum_{i=1}^n w_i(\mu)}, \quad w_i(\mu) = \exp\{-\frac{\gamma}{2}(x_i - \mu)^T V_0^{-1}(x_i - \mu)\} \quad (\text{繰返し重み付け平均})$$

(A3) もし  $\mu_t \in \{\mu_1, \dots, \mu_{t-1}\}$  とならばステップ(A2)を停止し、 $k=t-1$  と定め、そうでなければステップ(A2)を続ける。

**【簡単な例題】**



**【結論と課題】**

正規分布の平均を  $\gamma$  パワー・エントロピー最小化法による推定した。最尤推定値 ( $\gamma=0$ ) はデータがどうであれ、標本平均ベクトルで与えられるが、提案方法 ( $\gamma > 0$ ) では  $k$  平均を推定する。この意味で自発的にクラスター数  $k$  が決まる。提案方法は推定の強いロバスト性が示されていたが、このような大きな **モデル不確定性に対して柔軟な性能** が確認された。提案方法では分散行列は  $V_0$  として予め決められていた。クラスターによって共分散が異なる場合はマハラノビスの2乗距離:  $(\mu - \mu_j)^T V_j^{-1}(\mu - \mu_j)$  を使って(2)と同様にしてクラスターを決めることができる。ただし

$$V_j = \arg \min \{ L_\gamma(\mu_j, V) : V > O \} \quad (3)$$

今後の課題としてはアルゴリズムの計算量の節約、温度パラメータ  $\gamma$  のアニーリングについて考察が必要である。テンソルデータ、関数データ、変数選択 などについても今後の課題となる。