

多様体

Amari, Shun-ichi. "Information geometry of the EM and em algorithms for neural networks." Neural networks 8.9 (1995): 1379-1408.

交通・都市・国土学研究室

月田 光

前提

- ・ノイズの入った教師データから入出力のモデル化を行う際に、ニューラルネットワーク(NN)が用いられる。
- ・データが観測不能な値や欠損値(隠れ変数)を含む場合、隠れ変数を推定する。
→隠れ変数の推定に、EMアルゴリズムとemアルゴリズムが提案されている

EMアルゴリズム：条件付き期待値を用いた反復的な統計手法

emアルゴリズム：多様体においてKLダイバージェンスを反復的に最小化

目的

- ・EMアルゴリズムとemアルゴリズムがほとんどの場合等価であることを示す。
- ・EMとemに統一的な情報幾何学的枠組みを与え、その等価性を保証する条件を証明。
- ・例として以下を考える。
(1) 確率的多層パーセプトロン, (2) 混合エキスパートモデル, (3) 混合正規分布

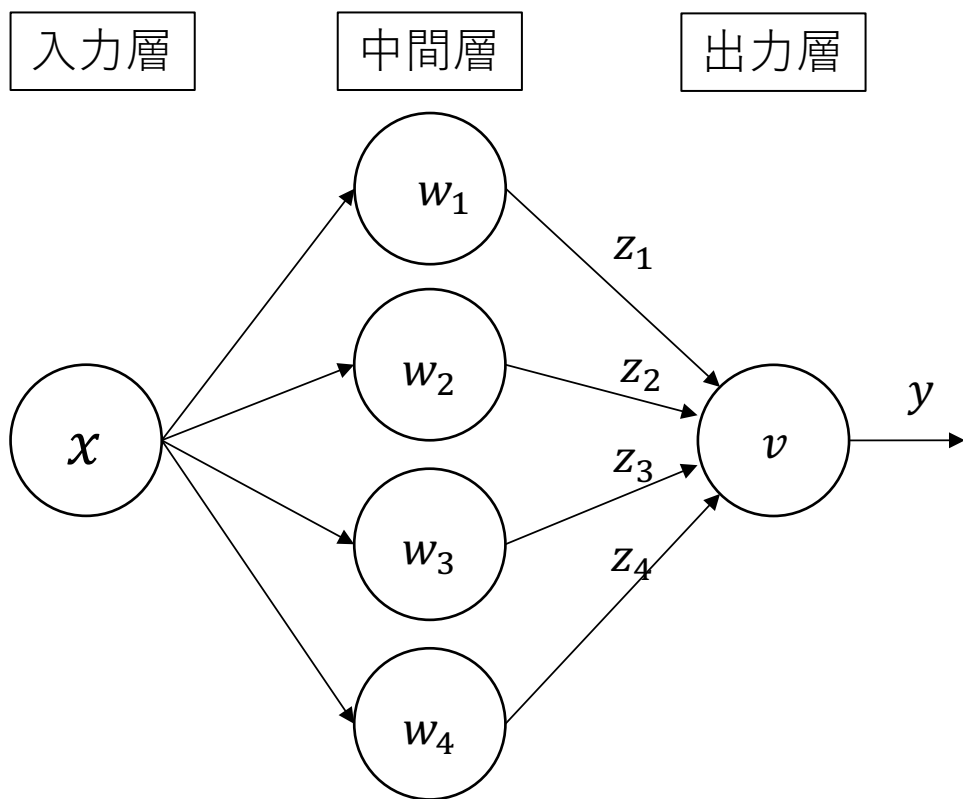
- 1.はじめに
- 2.指数型分布族とニューラルネットワーク
- 3.観測と推定の幾何学
- 4.隠れ変数と部分多様体
- 5.EMアルゴリズムとemアルゴリズム
- 6.EMとemの情報幾何学
- 7.拡張フレームワークにおけるEMとemの計算方法の等価性
- 8.学習の手順
- 9.emアルゴリズムの差分形式と増分形式
- 10.例 - 正規混合分布、放射状基底拡張、エキスパートネットの混合
- 11.結論

第1章はじめに 一部抜粋・和訳

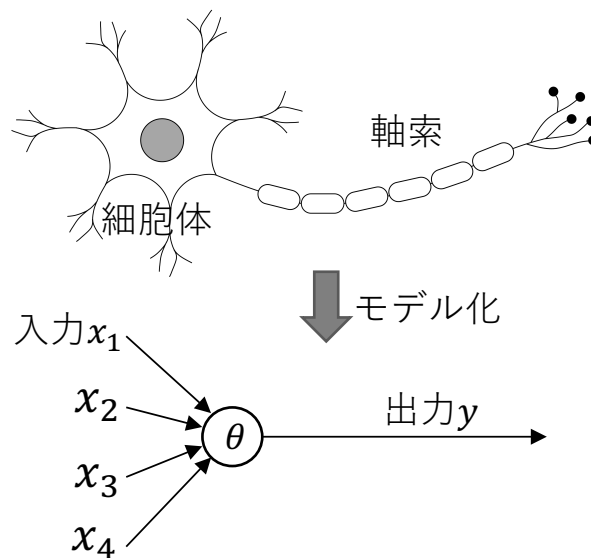
第2章は、基本的な統計モデルである指数族と曲線指数族についての紹介である。第3章は統計的予備知識で、最尤推定を幾何学的に説明する。ニューラルネットワークや統計学に慣れた読者はこれらのセクションを飛ばしてもよい。第4章では、本理論の主要な枠組みを紹介する。ニューラルネットワークの集合 M は、関連する確率分布の多様体 S に部分多様体として埋め込まれることが示される。問題は、 D と M の間のKullback-Leibler divergenceを最小にする、 M のネットワークと D のフィールドインデータを見つけることである。本論文の主要部分である第6章では、情報幾何学的な観点から両者の分析が行われる。第7章では、この2つのアルゴリズムがほとんどの実用的なケースで等価であることが示される。第8章では学習方法を扱い、第9章では M と D における双測地的勾配流を研究する。第10章では正規混合とエキスパートネットの混合を研究する。付録として、情報幾何学の直観的な入門を簡単に述べる。

ニューラルネットワーク(NN)

- ・ 入力を線形変換する計算処理が神経ネットワーク状に結合した数理モデル
- ・ 系全体で見ると、教師データから学習可能な非線形関数の普遍的な近似となる。
- ・ 入出力関係は x を入力したときの出力 y の条件付き確率 $p(y|x)$ で確率的に記述。
- ・ 機械学習のモデルとして用いられる。



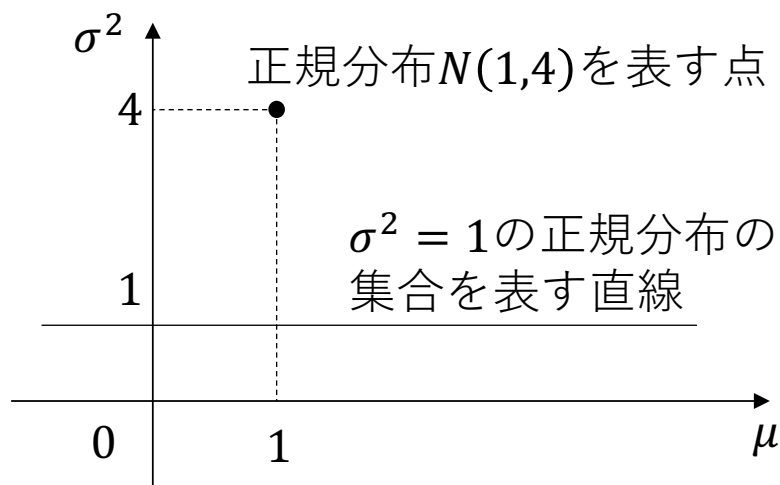
ベクトル $\theta = (\theta_1, \theta_2, \dots, \theta_n)$ で表されたパラメータを持つNNにおける確率分布 $P(y; \theta)$ や条件付き確率分布 $P(x|y; \theta)$ を考える。



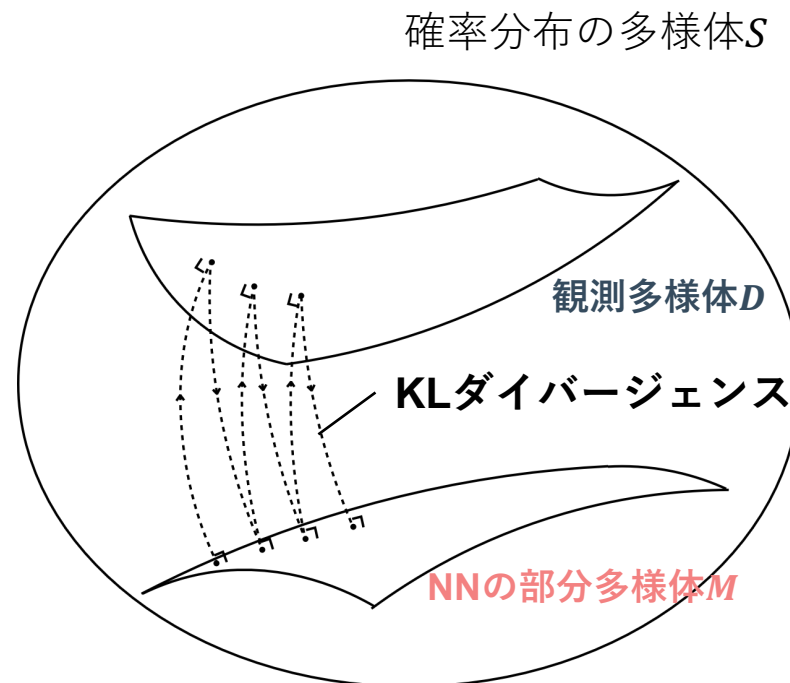
情報幾何学・確率分布の多様体

- 情報幾何は、確率分布の多様体の情報構造に由来し、新しい微分幾何学的概念を持つ新しい数学的枠組みとして発展してきた。

例：正規分布 $N(\mu, \sigma^2)$ を表す座標平面



パラメータの次元が $n = 2$ のため平面で、パラメータを一つ固定した部分集合は直線で表されているが、パラメータ数 n を一般化すると、これらは多様体、部分多様体となる。

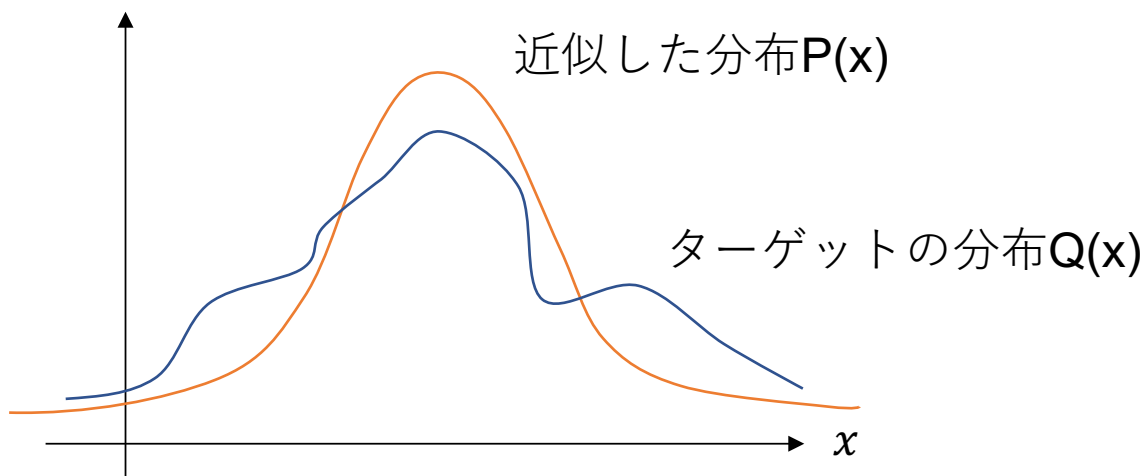


距離の尺度としてKLダイバージェンスを用い、これを最小化するM上の点を求める。

KLダイバージェンス(Kullback–Leibler divergence)

- 二つの確率分布の違いを数値化したもの.
- モデル化したある確率分布 $P(x)$ が, 観測された確率分布 $Q(x)$ をどれほど良く近似しているかを数値として表す.
- 非負の値をとり, ぴったり一致する場合に0となる.
- 計算方法は交差エントロピーから情報エントロピーを引く.

$$\begin{aligned}K(Q||P) &= H(Q, P) - H(Q) \\ &= E_{x \sim Q}[-\log P(x)] - E_{x \sim Q}[-\log Q(x)] \\ &= \int Q(x) \log \frac{Q(x)}{P(x)} dx\end{aligned}$$



指数型分布族の定義

- $r_i(x), i = 0, \dots, n$ を x の関数, ベクトル $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ をパラメータとし, 確率変数 x の確率密度関数が以下の式で表されるとき, その族 $S = \{p(x; \boldsymbol{\theta})\}$ を指数型分布族という. ただし, x はベクトルでも良い.

$$p(x; \boldsymbol{\theta}) = \exp \left\{ \sum_{i=1}^n \theta_i r_i(x) + k(x) - \psi(\boldsymbol{\theta}) \right\}$$

- x を \mathbf{r} に変換し, $\mathbf{r} = (r_1, \dots, r_n) \in R^n$ を R^n 上の適当な測度 $\mu(\mathbf{r})$ に関して以下の分布を持つベクトル確率変数として扱うことができる. ただし, 項 $k(x)$ とヤコビアンは $\mu(\mathbf{r})$ に吸収される.

$$p(\mathbf{r}; \boldsymbol{\theta}) = \exp \left\{ \sum_{i=1}^n \theta_i r_i(x) - \psi(\boldsymbol{\theta}) \right\}$$

例1: 正規分布

$$r_1 = x, r_2 = x^2, \theta_1 = \frac{\mu}{\rho^2}, \theta_2 = \frac{1}{2\rho^2}, \psi(\boldsymbol{\theta}) = \frac{\mu^2}{2\rho^2} + \log(\sqrt{2\pi}\rho)$$

例2: 離散分布 $p_i = \text{Prob}(x = i), i = 0, \dots, n$ とすると

$$r_i = \log \frac{p_i}{p_0}, \theta_i = \delta_i(x), \psi(\boldsymbol{\theta}) = -\log p_0 = \log\{1 + \sum \exp(\theta_i)\}$$

確率的多層パーセプトロン

右下の図のような出力が一つの隠れ層パーセプトロンを考える。

- ・ 入力を x , m 個隠れユニットの出力を $\mathbf{z} = (z_i), i = 1, \dots, m, z_i = 0 \text{ or } 1$ とする。
- ・ z_i の確率を次のように定義

$$p(z_i|x) = \varphi(z_i, \mathbf{w}_i \cdot \mathbf{x})$$

ただし、 φ はシグモイド関数

$$\varphi(z, u) = \frac{\exp(zu)}{1 + \exp(u)}$$

- ・ \mathbf{w}_i は i 番目の隠れユニットのシナプスの重みベクトル
- ・ \mathbf{x} に $x_0 \equiv 1$ を追加し, $\mathbf{w}_i \cdot \mathbf{x} = \sum_{j=0}^n w_{ij} x_j$ となる。
- ・ $p(z_i = 0|x) + p(z_i = 1|x) = 1$ になる。

次に出力ユニットを考える。

- ・ 隠れ層から信号 \mathbf{z} を受け取り, 2 値出力 y を返す。
- ・ y の確率は \mathbf{z} に依存し, その確率を次のようにおく。

$$p(y|\mathbf{z}) = \varphi(y, \mathbf{v} \cdot \mathbf{z})$$

ただし, \mathbf{v} は出力ユニットのシナプスの重みベクトル \mathbf{z} に $z_0 \equiv 1$ を追加し, 定数項も考慮。

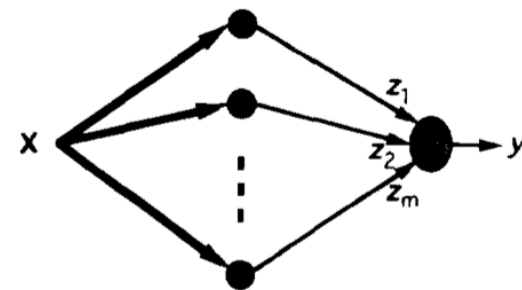


FIGURE 2. Stochastic multilayer perceptron.

確率的多層パーセプトロン

- ・ 入力 x が与えられた時, 変数 (y, z) の条件付き確率は以下の式となる.

$$p(y, z|x; \mathbf{u}) = p(y|z; \mathbf{u}) \prod_{i=1}^m p(z_i|x; \mathbf{u})$$

- ・ ただし, \mathbf{u} は全パラメータを要約した変数である.

$$\mathbf{u} = (\mathbf{w}_1, \dots, \mathbf{w}_m, \mathbf{v})$$

- ・ 2値出力のパーセプトロンの期待値や, $[0,1]$ の出力をする決定論的パーセプトロンでは以下の値となる.

$$z_i = \varphi(\mathbf{1}, \mathbf{w}_i \cdot x), y = (1, \mathbf{v} \cdot \mathbf{z})$$

- ・ x が確率分布 $q(x) > 0$ に従って生成されるとき, (y, z, x) の全確率分布は次の式.

$$p(y, z, x; \mathbf{u}) = q(x)p(y, z|x; \mathbf{u})$$

- ・ 対数を取ると次の式になる.

$$l(y, z, x; \mathbf{u}) = \log p(y, z|x; \mathbf{u}) + \log q(x)$$

- ・ x の分布 $q(x)$ が未知でも, 条件付き確率の対数 $\log p(y, z|x; \mathbf{u})$ を最大化することが $l(y, z, x; \mathbf{u})$ の最大化である.

確率的多層パーセプトロンと指数型分布族

- ・対数をさらに変形する.

$$l(y, \mathbf{z}, \mathbf{x} | \mathbf{u}) = y\mathbf{v} \cdot \mathbf{z} + \sum_i z_i \mathbf{w}_i \cdot \mathbf{x} - \log\{1 + \exp \mathbf{v} \cdot \mathbf{z}\} - \sum_i \log\{1 + \exp(\mathbf{w}_i \cdot \mathbf{x})\}$$

- ・ \mathbf{z} の値として、各成分が $\{0,1\}$ をとる m 次元ベクトル \mathbf{k} を導入. またデルタ関数を導入.
 $\delta_{\mathbf{k}}(\mathbf{z})$ は \mathbf{k} でインデックスされた 2^m 次元のベクトル確率変数

$$\delta_{\mathbf{k}}(\mathbf{z}) = \begin{cases} 1 & \mathbf{z} = \mathbf{k} \\ 0 & \mathbf{z} \neq \mathbf{k} \end{cases}$$

- ・また $\mathbf{w}_{\mathbf{k}} = \sum_i \mathbf{k}_i \mathbf{w}_i$ を置き, $\sum_i z_i \mathbf{w}_i \cdot \mathbf{x} = \sum_{\mathbf{k}} \delta_{\mathbf{k}}(\mathbf{z}) \mathbf{w}_{\mathbf{k}} \cdot \mathbf{x}$ とすることで, 次式となる.

$$\begin{aligned} l(y, \mathbf{z} | \mathbf{x}, \mathbf{u}) &= \sum_{\mathbf{k}} \mathbf{k} \cdot \mathbf{v} y \delta_{\mathbf{k}}(\mathbf{z}) - \sum_{\mathbf{k}} \delta_{\mathbf{k}}(\mathbf{z}) \log\{1 + \exp(\mathbf{k} \cdot \mathbf{v})\} \\ &\quad + \sum_{\mathbf{k}} \delta_{\mathbf{k}}(\mathbf{z}) \mathbf{w}_{\mathbf{k}} \cdot \mathbf{x} - \sum_i \log\{1 + \exp(\mathbf{w}_i \cdot \mathbf{x})\} \end{aligned}$$

確率的多層パーセプトロンと指数型分布族

- ・ 以下のようにおくと, $p(y, \mathbf{z} | \mathbf{x}, \mathbf{u})$ が指数型分布族であるとわかる.

$$r_{1,k} = y\delta_k(\mathbf{z}), r_{2,k} = \delta_k(\mathbf{z}), \theta_{1,k} = \mathbf{k} \cdot \mathbf{v}, \theta_{2,k} = \mathbf{w}_k \cdot \mathbf{x} - \log\{1 + \exp(\mathbf{k} \cdot \mathbf{v})\}$$

$$\psi(\boldsymbol{\theta}) = \sum_i \log\{1 + (\mathbf{w}_i \cdot \mathbf{x})\}$$

$$p(\mathbf{r} | \mathbf{x}; \mathbf{u}) = \exp\{\mathbf{r} \cdot \boldsymbol{\theta} - \psi(\boldsymbol{\theta})\}$$

繰り返し観測と積空間

- $\mathbf{r}_1, \dots, \mathbf{r}_T$ を同じ分布 $p(\mathbf{r}; \boldsymbol{\theta})$ からの T 個の独立な観測値とする.

$$p(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_T; \boldsymbol{\theta}) = \prod_{t=1}^T p(\mathbf{r}_t; \boldsymbol{\theta})$$

- $p(\mathbf{r}; \boldsymbol{\theta})$ が指数型の場合

$$p(\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_T; \boldsymbol{\theta}) = \exp\{(\sum \mathbf{r}_t) \cdot \boldsymbol{\theta} - T\psi(\boldsymbol{\theta})\}$$

- よって全体としても指数型分布族になる.
- 次に条件付き分布 $p(\mathbf{y}, \mathbf{z} | \mathbf{x})$ を考える.
- $(\mathbf{x}_t, \mathbf{y}_t, \mathbf{z}_t), t = 1, \dots, T$ を指数型分布族からの独立した観測値とする.
- $p\{\mathbf{r}_t; \boldsymbol{\theta}(\mathbf{x}_t)\} = \exp\{\boldsymbol{\theta}(\mathbf{x}_t) \cdot \mathbf{r}_t - \psi_t\}$ とすると, 結合条件分布は次の式.

$$p\{(\mathbf{y}_1, \mathbf{z}_1), \dots, (\mathbf{y}_T, \mathbf{z}_T) | \mathbf{x}_1, \dots, \mathbf{x}_T\} = \exp\left\{\sum \boldsymbol{\theta}(\mathbf{x}_t) \cdot \mathbf{r}_t - \sum \psi_t\right\}$$

- ただし, $\mathbf{r} = \mathbf{r}(\mathbf{y}, \mathbf{z})$

繰り返し観測と積空間

- $S_t = \{p(r_t; \theta(x_t))\}$ は t 番目の入力に対する条件付き確率分布の多様体であり、全部の結合分布は以下の積空間に与えられる。

$$S_T^* = S_1 \times S_2 \times \cdots \times S_T$$

- S_T^* の観測値 \mathbf{r}_T^* を次のようにおく。

$$\mathbf{r}_T^* = (\mathbf{r}_1, \dots, \mathbf{r}_T)$$

- パラメータ $\boldsymbol{\theta}_T^*$ を次のようにおく。

$$\boldsymbol{\theta}_T^* = \{\boldsymbol{\theta}(x_1), \dots, \boldsymbol{\theta}(x_T)\}$$

- $\boldsymbol{\theta}_T^* \cdot \mathbf{r}_T^* = \sum_{t=1}^T \boldsymbol{\theta}_t \cdot \mathbf{r}_t$, $\psi(\boldsymbol{\theta}_T^*) = \sum_{t=1}^T \psi(\boldsymbol{\theta}_t)$ とおくと次式となり、指数型分布族。

$$p(\mathbf{r}_T^*; \boldsymbol{\theta}_T^*) = \exp\{\boldsymbol{\theta}_T^* \cdot \mathbf{r}_T^* - \psi(\boldsymbol{\theta}_T^*)\}$$

- ただし、 $\boldsymbol{\theta}_T^*$ は $\theta_{1,\mathbf{k}} = \mathbf{k} \cdot \mathbf{v}$, $\theta_{2,\mathbf{k}} = \mathbf{w}_{\mathbf{k}} \cdot \mathbf{x} - \log\{1 + \exp(\mathbf{k} \cdot \mathbf{v})\}$ より \mathbf{w}_i と \mathbf{v} に制限され、任意の値をとることができず、 S_T^* の部分領域に制限される

曲線指数型分布族

- $S = \{p(\mathbf{r}; \boldsymbol{\theta})\}$ を n 次元指数型分布族とし、 M をその m 次元部分多様体とすると、 S の一部を占める部分集合 M は曲線指数型分布族(Curved Exponential Families) という。 M は S の部分多様体である。
- $\mathbf{u} = (u_1, \dots, u_m)$, $m \leq n$ を M の座標系とすると、 M は \mathbf{u} の関数 $\boldsymbol{\theta}(\mathbf{u})$ を用いて次の分布により構成される

$$p(\mathbf{r}; \boldsymbol{\theta}(\mathbf{u})) = \exp\{\boldsymbol{\theta}(\mathbf{u}) \cdot \mathbf{r} - \psi(\boldsymbol{\theta}(\mathbf{u}))\}$$

- $\boldsymbol{\theta}_T^* = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T\}$ は積空間 S_T^* の座標であり、パラメータ \mathbf{u} を用いて $\boldsymbol{\theta}_t = \boldsymbol{\theta}(\mathbf{x}_t, \mathbf{u})$.
- \mathbf{x}_t が与えられた時、 $\boldsymbol{\theta}_T^*$ における自由なパラメータは \mathbf{u} であり、 T が増加しても次元数は固定。
- よって、 M_T^* は S_T^* の部分多様体を形成し、 M 上の分布 $p(\mathbf{r}; \boldsymbol{\theta}(\mathbf{u}))$ は S に属しており、 $\boldsymbol{\theta}(\mathbf{u})$ は S 上での座標である(右図)。

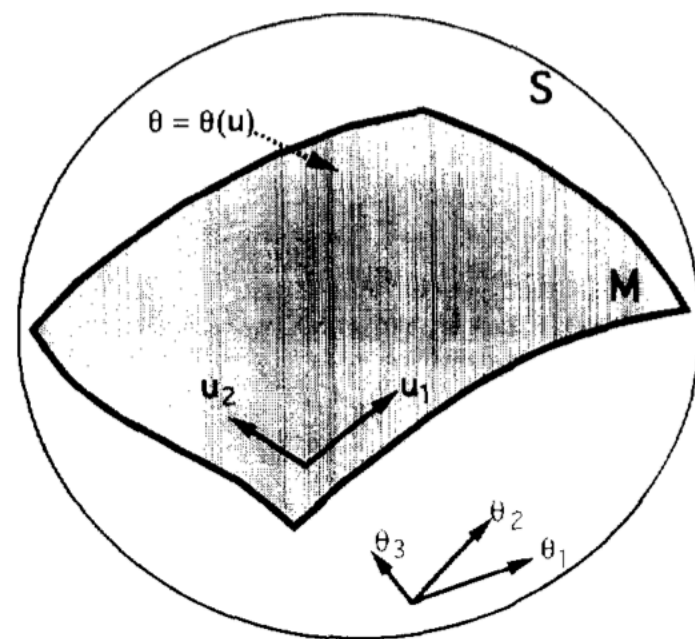


FIGURE 3. Curved exponential family.

曲線指数型分布族の例：正規分布の乗算モデル

- ε を正規分布 $N(0,1)$ に従う確率変数とする.
- ノイズ ε が入った大きさ1の信号を入力, u 倍に増幅された出力 x がされるとする.
$$x = u(1 + \varepsilon)$$
- x は $N(u, u^2)$ に従う.
- この集合 M は正規分布 $S = \{N(\mu, \rho^2)\}$ の曲線指数型分布族である. 平均 μ と分散 ρ^2 は独立ではなく, 共通の u によって規定される.
- よって, M は S 上で曲線(放物線)となる

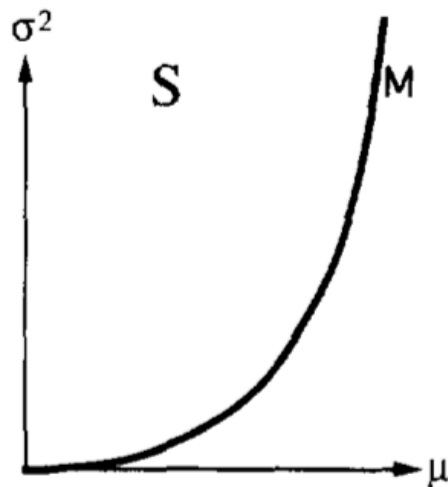


FIGURE 4. Example of curved exponential family.

期待値パラメータ

- すべての確率変数 \mathbf{r} が可視であると仮定する.
- パラメータ $\boldsymbol{\theta}$ で表される分布 $p(\mathbf{r}; \boldsymbol{\theta})$ の期待値を $\boldsymbol{\eta}$ とする.

$$\boldsymbol{\eta}(\boldsymbol{\theta}) = E_{\boldsymbol{\theta}}[\mathbf{r}] = \int \mathbf{r} p(\mathbf{r}; \boldsymbol{\theta}) d\mu(\mathbf{r})$$

- ただし, 確率の積分 $\int p(\mathbf{r}; \boldsymbol{\theta}) d\mu(\mathbf{r}) = \int \exp\{\boldsymbol{\theta} \cdot \mathbf{r} - \psi(\boldsymbol{\theta})\} d\mu(\mathbf{r})$ は1であることより, これを $\boldsymbol{\theta}$ で微分することで, 以下の式が導ける.

$$\boldsymbol{\eta}(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} \psi(\boldsymbol{\theta})$$

- さらに, $\boldsymbol{\theta}$ と $\boldsymbol{\eta}$ の変換は1対1であるため, $\boldsymbol{\eta}$ は S のもう一つのペラメータとして使用できる. ここで, $\boldsymbol{\eta}$ を期待値パラメータと呼ぶ.
- よって, $\boldsymbol{\theta} = \boldsymbol{\theta}(\boldsymbol{\eta})$ と逆変換できる.

- $\boldsymbol{\eta}$ で定義される確率分布 $p^*\{\mathbf{r}; \boldsymbol{\eta}(\boldsymbol{\theta})\} = p(\mathbf{r}; \boldsymbol{\theta})$ を導入.
- $\varphi(\boldsymbol{\eta})$ を分布 $p^*\{\mathbf{r}; \boldsymbol{\eta}(\boldsymbol{\theta})\}$ のエントロピーの負とする.

$$\varphi(\boldsymbol{\eta}) = \int p^*(\mathbf{r}; \boldsymbol{\eta}) \log p^*(\mathbf{r}; \boldsymbol{\eta}) d\mu(\mathbf{r})$$

- すると, $\boldsymbol{\theta}(\boldsymbol{\eta})$ は次のように与えられる (Amari, 1985).

$$\boldsymbol{\theta}(\boldsymbol{\eta}) = \frac{\partial}{\partial \boldsymbol{\eta}} \varphi(\boldsymbol{\eta})$$

期待値パラメータの例

- ・ 正規分布の場合

$$\eta_1 = E[x] = \mu = \frac{\theta_1}{2\theta_2}, \eta_2 = E[x^2] = \mu^2 + \rho^2 = \left(\frac{\theta_1}{2\theta_2}\right)^2 - \frac{1}{2\theta_2}$$

- ・ 離散分布の場合

$$\eta_i = E[\delta_i(x)] = p_i = \frac{\exp(\theta_i)}{1 + \sum \exp(\theta_i)}$$

- ・ 確率的パーセプトロン $S = \{p(y, \mathbf{z}; \boldsymbol{\theta})\}$ の場合

$$\boldsymbol{\eta} = (\eta_{1,k}, \eta_{2,k})$$

$$\eta_{1,k} = E[y\delta_k(\mathbf{z})] = \varphi(1, \mathbf{v} \cdot \mathbf{k}) \prod \varphi(k_i, \mathbf{w}_i \cdot \mathbf{x})$$

$$\eta_{2,k} = E[\delta_k(\mathbf{z})] = \prod \varphi(k_i, \mathbf{w}_i \cdot \mathbf{x})$$

期待値パラメータと最尤推定量

- ・ 期待値パラメータは最尤推定量の研究に有用である.
- ・ 指数型分布族 $S = \{p(\mathbf{r}; \boldsymbol{\theta})\}$ の尤度関数の対数は以下の式.

$$l(\check{\mathbf{r}}; \boldsymbol{\theta}) = \boldsymbol{\theta} \cdot \check{\mathbf{r}} - \psi(\boldsymbol{\theta})$$

- ・ これを微分することにより最尤推定量が求められる.

$$\frac{\partial l}{\partial \boldsymbol{\theta}} = \check{\mathbf{r}} - \boldsymbol{\eta}(\hat{\boldsymbol{\theta}}) = \mathbf{0}$$

- ・ よって期待値 $\boldsymbol{\eta}$ は観測値そのもの. $\hat{\boldsymbol{\theta}}$ も $\hat{\boldsymbol{\theta}} = \frac{\partial}{\partial \boldsymbol{\eta}} \varphi(\hat{\boldsymbol{\eta}})$ から求められる.
- ・ $\hat{\boldsymbol{\eta}}$ または $\hat{\boldsymbol{\theta}}$ は S における点を決定するものであり, これを観測点と呼ぶ.

$\check{\mathbf{r}}$: \mathbf{r} の観測値
 $\hat{\boldsymbol{\theta}}$: $\boldsymbol{\theta}$ の最尤推定量

繰り返し観測と観測点

- $\mathbf{r}_1, \dots, \mathbf{r}_T$ を T 個の独立な確率変数とし, \mathbf{r}_t は分布 $p(\mathbf{r}_t; \boldsymbol{\theta}_t)$ に従うとする.

$$p(\mathbf{r}_1, \dots, \mathbf{r}_T; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_T) = \prod_{t=1}^T p(\mathbf{r}_t; \boldsymbol{\theta}_t)$$

$$p(\mathbf{r}_T^*; \boldsymbol{\theta}_T^*) = \exp\{\boldsymbol{\theta}_T^* \cdot \mathbf{r}_T^* - \psi\}$$

- 期待座標 $\boldsymbol{\eta}_T^* = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_T)$ は $\boldsymbol{\eta}_t = E_{\boldsymbol{\theta}_t}[\mathbf{r}_t]$ で与えられ, 観測値 $\check{\mathbf{r}}_t$ により以下の通り.

$$\boldsymbol{\eta}_T^* = (\check{\mathbf{r}}_1, \dots, \check{\mathbf{r}}_T)$$

- 対応する最尤推定量 $\hat{\boldsymbol{\theta}}_t$ は次のように与えられる.

$$\hat{\boldsymbol{\theta}}_t = \frac{\partial \varphi(\hat{\boldsymbol{\eta}}_t)}{\partial \boldsymbol{\eta}}$$

- NN の場合, $\boldsymbol{\theta}_t, \boldsymbol{\eta}_t$ は入力 \mathbf{x}_t とパラメータ \mathbf{u} により $\boldsymbol{\theta}_t = \boldsymbol{\theta}(\mathbf{x}_t, \mathbf{u}), \boldsymbol{\eta}_t = \boldsymbol{\eta}(\mathbf{x}_t, \mathbf{u})$ と書かれ, これが曲線指数型分布族 M^* を定義する. しかし観測点 $\hat{\boldsymbol{\eta}}_T^*, \hat{\boldsymbol{\theta}}_T^*$ がこれを満たすとは限らず, 必ずしも M に含まれるとは限らない.
- よって最尤推定量 $\hat{\mathbf{u}}$ を得るためには, M に射影する必要がある.

繰り返し観測と観測点

- すべての \mathbf{r}_t が同じ分布 $p(\mathbf{r}_t; \boldsymbol{\theta})$ に従うとする. $\boldsymbol{\theta} = \boldsymbol{\theta}_1 = \dots = \boldsymbol{\theta}_T$

$$p(\mathbf{r}_1, \dots, \mathbf{r}_T; \boldsymbol{\theta}) = \exp \left\{ \sum_{t=1}^T \mathbf{r}_t \cdot \boldsymbol{\theta} - T\psi(\boldsymbol{\theta}) \right\}$$

- 観測値の平均値 $\bar{\mathbf{r}} = \frac{1}{T} \sum_{t=1}^T \mathbf{r}_t$ を導入し, S_T^* を単一の S に還元できる.

$$p(\bar{\mathbf{r}}; \boldsymbol{\theta}) = \exp \{ T(\bar{\mathbf{r}} \cdot \boldsymbol{\theta}) + Tk(\bar{\mathbf{r}}) - T\psi(\boldsymbol{\theta}) \}$$

- これは $\bar{\mathbf{r}}$ が $\boldsymbol{\theta}$ や \mathbf{u} を推定するのに十分な統計量であることを示している.
- ただし, $k(\bar{\mathbf{r}})$ は \mathbf{r} の変換によって生じ, $\exp\{Tk(\bar{\mathbf{r}})\} = \int_{\bar{\mathbf{r}}} \prod d\mu(\mathbf{r}_t)$ で表され, $d\bar{\mu}(\bar{\mathbf{r}}) = \exp\{Tk(\bar{\mathbf{r}})\}d\mu$ で除去できる.
- i.i.d.である場合は, S_T^* を参照せず指数型関数族 S で $\bar{\mathbf{r}} \cdot \boldsymbol{\theta} - \psi(\boldsymbol{\theta})$ を最大化することで最尤推定量 $\hat{\boldsymbol{\theta}}$ が与えられる.
- 以上より, $\bar{\mathbf{r}} = \frac{\partial}{\partial \boldsymbol{\theta}} \psi(\boldsymbol{\theta}), \hat{\boldsymbol{\eta}} = \bar{\mathbf{r}}$. 観測されたデータ $\check{\mathbf{r}}_t$ はFisherの意味で情報損失無しに S における $\boldsymbol{\eta}$ 座標系の一点 $\hat{\boldsymbol{\eta}}$ に要約される.

曲線指数型分布族での推定

- S 上の M における最尤推定を考える。
- 観測点 $\hat{\eta} = \bar{r}$ が与えられた時、最尤推定量 \hat{u} は次式対数尤度最大化で求める。

$$l(\bar{r}; \mathbf{u}) = \bar{r} \cdot \boldsymbol{\theta}(\mathbf{u}) - \psi\{\boldsymbol{\theta}(\mathbf{u})\}$$

- これは、 $p(\bar{r}; \hat{\boldsymbol{\theta}})$ から $p(\bar{r}; \boldsymbol{\theta}(\mathbf{u}))$ へのKLダイバージェンスを最小化することと等価。

$$\begin{aligned} K(\hat{\boldsymbol{\theta}} || \boldsymbol{\theta}(\mathbf{u})) &= \int p(\bar{r}; \hat{\boldsymbol{\theta}}) \log \frac{p(\bar{r}; \hat{\boldsymbol{\theta}})}{p(\bar{r}; \boldsymbol{\theta}(\mathbf{u}))} d\mu(\bar{r}) \\ &= -H(\hat{\boldsymbol{\theta}}) - \hat{\eta} \cdot \boldsymbol{\theta}(\mathbf{u}) + \psi(\boldsymbol{\theta}(\mathbf{u})) \end{aligned}$$

- 幾何学的には $\hat{\boldsymbol{\theta}}$ を M に射影することで与えられる。
- 尤度方程式を解くことで \mathbf{u} を求める。

$$\frac{\partial}{\partial \mathbf{u}} \{\bar{r} \cdot \boldsymbol{\theta}(\mathbf{u}) - \psi(\boldsymbol{\theta}(\mathbf{u}))\} = \mathbf{0}$$

$$\mathbf{B}(\mathbf{u})\{\bar{r} - \boldsymbol{\eta}(\mathbf{u})\} = \mathbf{0}$$

$$\mathbf{B}(\mathbf{u}) = \frac{\partial \boldsymbol{\theta}(\mathbf{u})}{\partial \mathbf{u}} = \left(\frac{\partial \theta_i(\mathbf{u})}{\partial u_j} \right)$$

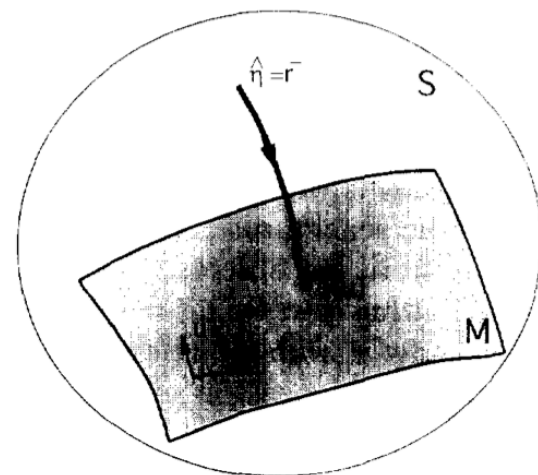


FIGURE 5. Maximum likelihood estimation.

曲線指数型分布族での推定

- 直積空間 S_T^* で KL ダイバージェンスは

$$K(\hat{\boldsymbol{\theta}}_T^* \parallel \boldsymbol{\theta}_T^*(\mathbf{u})) = \sum_{t=1}^T K(\hat{\boldsymbol{\theta}}_t \parallel \boldsymbol{\theta}_t(\mathbf{u})) = - \sum H(\hat{\boldsymbol{\theta}}_t) - \sum \hat{\boldsymbol{\eta}}_t \cdot \boldsymbol{\theta}_t(\mathbf{u}) + \sum \psi(\boldsymbol{\theta}_t(\mathbf{u}))$$

- 最尤推定量 $\hat{\mathbf{u}}$ は M^* における $\boldsymbol{\theta}_T^*(\mathbf{u})$ のうち、観測点 $\hat{\boldsymbol{\eta}}_T^* = \bar{\mathbf{r}}^* \in S_T^*$ に最も近い点。
- 尤度式 $\sum_{t=1}^T \mathbf{B}(\mathbf{x}_t, \mathbf{u}) \{\check{\mathbf{r}}_t - \boldsymbol{\eta}(\mathbf{x}_t, \mathbf{u})\} = \mathbf{0}$

ただし、 $\mathbf{B}(\mathbf{x}_t, \mathbf{u}) = \frac{\partial \boldsymbol{\theta}(\mathbf{x}_t, \mathbf{u})}{\partial \mathbf{u}}$

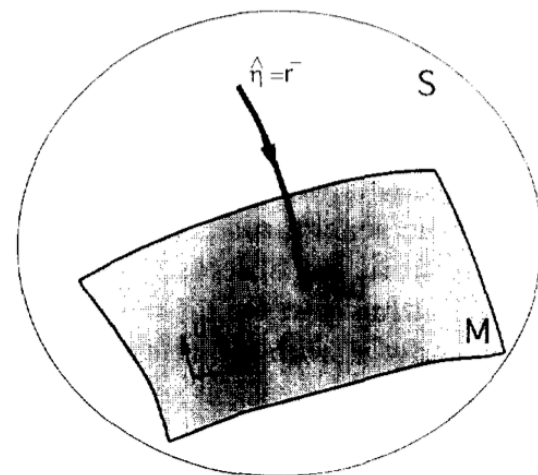


FIGURE 5. Maximum likelihood estimation.

部分観測

- 確率変数 \mathbf{r} の一部が観測できない場合を考える。
例) NNの隠れ層の出力は観測されない。
- この場合、観測されない変数も適切に決定する必要がある。

$\mathbf{r} = (\mathbf{s}_v, \mathbf{s}_h)$ とおく。
 \mathbf{s}_v : 観測可能な統計値
 \mathbf{s}_h : 観測不可な欠損値

- この時、観測点 $\hat{\boldsymbol{\eta}} = \mathbf{r} = (\mathbf{s}_v, \mathbf{s}_h)$ を一意に特定することができず、任意の \mathbf{s}_h をとる。
- この候補点は観測データの部分多様体 D を形成する。

$$D = \{\hat{\boldsymbol{\eta}} \mid \hat{\boldsymbol{\eta}} = \mathbf{r}(\mathbf{s}_v, \mathbf{s}_h), \mathbf{s}_v \text{は観測値(固定)}, \mathbf{s}_h \text{は任意}\}$$

- 多くの場合、 \mathbf{r} は \mathbf{s}_h に対して線形
 - $\mathbf{r}(\mathbf{s}_v, \mathbf{s}_{1h}), \mathbf{r}(\mathbf{s}_v, \mathbf{s}_{2h}) \in D$, $0 < t < 1$ に対し
 $t\mathbf{r}(\mathbf{s}_v, \mathbf{s}_{1h}) + (1-t)\mathbf{r}(\mathbf{s}_v, \mathbf{s}_{2h}) \in D$
- h の次元の部分多様体を形成し、データ部分多様体、観測部分多様体と呼ばれる。
- D が離散ベクトル変数の時、 D は多数の候補点からなるが、 $\boldsymbol{\eta}$ 座標の線形部分多様体に拡張する。

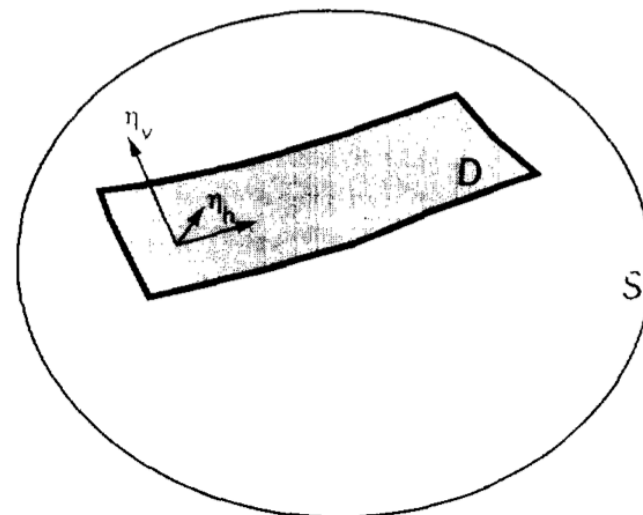
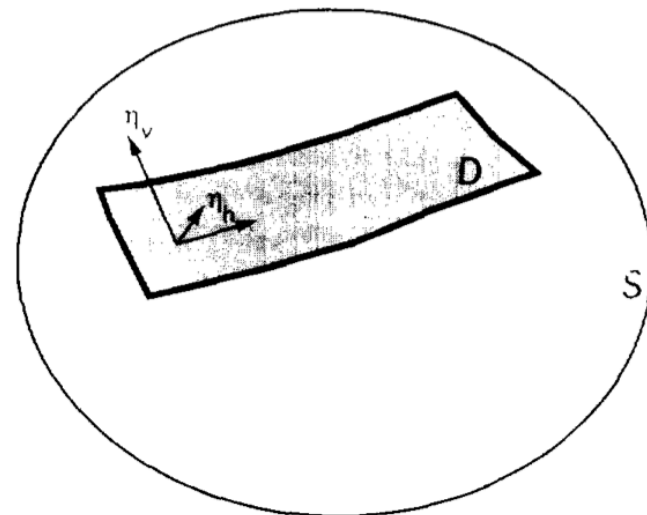


FIGURE 6. Data submanifold D : $\boldsymbol{\eta} = \text{fixed}$.

部分観測

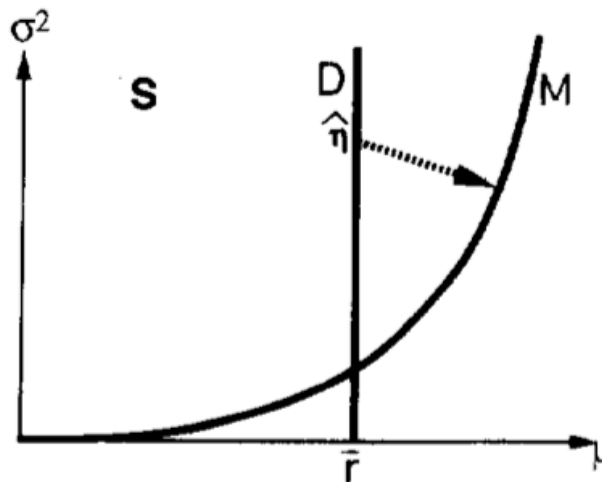
- η_h は D の内部座標系の役割を果たす.
- D が η 座標で線形である場合, 非特異行列 A を用いて $\tilde{\eta} = A\eta, \tilde{\mathbf{r}} = A\mathbf{r}$ と線形変換し $\tilde{\eta} = (\tilde{\eta}_v, \tilde{\eta}_h)$ に分割できる. またこの時 $\tilde{\boldsymbol{\theta}} = A^{-1}\boldsymbol{\theta}$ であり, 確率分布は以下の通り

$$p(\tilde{\mathbf{r}}; \tilde{\boldsymbol{\theta}}) = \exp\{\tilde{\boldsymbol{\theta}} \cdot \tilde{\mathbf{r}} - \tilde{\psi}(\tilde{\boldsymbol{\theta}})\}$$
- $\tilde{\boldsymbol{\theta}}$ が新しい正規化パラメータとなり, $\tilde{\boldsymbol{\eta}} = \frac{\partial \tilde{\psi}}{\partial \tilde{\boldsymbol{\theta}}}$ が期待値パラメータとなる.

FIGURE 6. Data submanifold D : η = fixed.

例：正規分布の乗法モデルの部分観測

- $x_1, \dots, x_T : N(u, u^2)$ に従う確率変数.
- $r_v = \frac{1}{T} \sum x_t$ (観測可能), $r_h = \frac{1}{T} \sum x_t^2$ (観測不可)
- モデル多様体 M は曲線指数型分布族. η 座標は $\eta_v = \mu = u, \eta_h = \mu^2 + \sigma^2 = 2u^2$
- データ多様体 D は $D = \{N(\mu, \sigma^2) | \mu = \tilde{r}, \sigma \text{ は任意}\}$
- $\sum x_t$ と $\sum x_t^2$ が観測される時観測点は一意的だが, この場合正確な観測点は D 上にあることしかわからない.
- 最尤推定量は $\hat{\eta}$ を m 射影することで得られる,

FIGURE 7. Example of data submanifold D and observed point $\hat{\eta}$.

例：確率的パーセプトロンのデータ部分多様体

- 時刻 t の確率変数 $\mathbf{r}(t) = (\mathbf{r}_{1,k,t}, \mathbf{r}_{2,k,t})$ において、観測点 $\hat{\boldsymbol{\eta}}$ が次のように表される。

$$\hat{\boldsymbol{\eta}}_{1,k} = \mathbf{r}_{1,k,t} = y_t \delta_{\mathbf{k}}(\mathbf{z}_t), \quad \hat{\boldsymbol{\eta}}_{2,k} = \mathbf{r}_{2,k,t} = \delta_{\mathbf{k}}(\mathbf{z}_t)$$

- y_t が観測可能、 $\delta_{\mathbf{k}}(\mathbf{z}_t)$ が観測不可能であり、 \mathbf{z}_t により候補点が存在。
- η 座標で線形結合をとり、 $\delta_{\mathbf{k}}(\mathbf{z}_t)$ につけた重みを $\alpha_{\mathbf{k}}$ とおくと、観測多様体 D_t は

$$D_t = \{\eta \mid \eta_{1,k} = y_t \alpha_{\mathbf{k}}, \eta_{2,k} = \alpha_{\mathbf{k}}\}$$

ただし、 $\sum_{\mathbf{k}} \alpha_{\mathbf{k}} = 1$

- $\alpha_{\mathbf{k}}$ は $\mathbf{z}_t = \mathbf{k}$ の確率と解釈しても良いが、観測不可能。
- データ部分多様体は

$$D_T^* = D_1 \times \cdots \times D_T$$

EM(Expectation and Maximizing)アルゴリズム

- ・部分的に観測されたデータ D または D_T^* から良い推定量 $\hat{\mathbf{u}}$ を得る手法.
- ・真のパラメータ \mathbf{u} と欠損データを同時に反復的に推定する.
- ・ \mathbf{u} に依存する分布の候補があれば $\hat{\mathbf{r}}$ を推測することができ, $\hat{\mathbf{r}}$ があれば \mathbf{u} が推測可能

方法

0. 準備. 任意の初期推測値 $\hat{\mathbf{u}}_0$ を選ぶ. 初期推測分布 $\hat{P}_0 \in M$ は $\theta(\hat{\mathbf{u}}_0)$ で与えられる.

1. Eステップ. 候補確率分布 $\hat{P}_i \in M$ に基づき, \mathbf{r} の条件付き期待値を計算する.

$$\hat{\boldsymbol{\eta}}_{(i)} = E[\mathbf{r} | \mathbf{s}_v; \hat{P}_i]$$

これが座標 $\hat{\boldsymbol{\eta}}_{(i)}$ で表される観測点 $\hat{Q}_i \in D$ の i 番目の候補となる.

2. Mステップ. 推定された観測点 \hat{Q}_i から最尤推定量 $\hat{\mathbf{u}}_{i+1}$ を算出, パラメータを更新.
これはKLダイバージェンス $K(\hat{Q}_i || P)$ の最小化または対数尤度の最大化で求める.

パラメータが収束するまでEステップとMステップを繰り返す.

EM(Expectation and Maximizing)アルゴリズム

- 部分観測データが \mathbf{r} でまとめられない場合，積空間 S_T^* にデータ多様体 D_T^* が存在。この時，Eステップは次のように推定する。

$$\hat{\mathbf{r}}_T^* = E[\mathbf{r}_T^* | s_{v,1}, \dots, s_{v,T}; \hat{P}_i]$$

ただし， \mathbf{r}_t は $s_{v,t}$ のみと相関があるため， $\hat{\eta}_t(i) = E[\mathbf{r}_t | s_{v,t}; \hat{P}_i]$

- 尤度は1回のEステップとMステップで増加する。
→尤度関数の局所最大に収束する(最尤推定量に収束するとは限らない)。
- 対数尤度を直接最大化する勾配法に比べ，大域的収束特性を持つ。

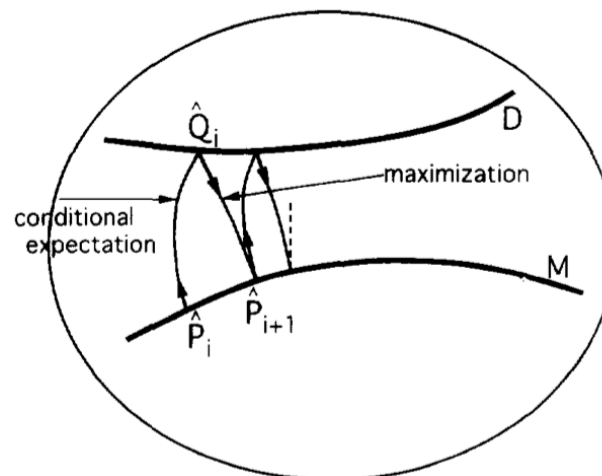


FIGURE 8. EM algorithm.

em(exponential mixture)アルゴリズム

- モデルは多様体 M に、観測値は多様体 D にあるため、 D と M のダイバージェンスを最小にする $\hat{P} \in M, \hat{Q} \in D$ の組を探索する。 \hat{Q} は最尤推定量となる。

$$K(\hat{Q}||\hat{P}) = \min_{P \in M, Q \in D} K(Q||P)$$

- 与えられた Q に対し $K(Q||P)$ を最小化する $\hat{P} \in M$ は Q の M への m -射影で与えられる。
 - m -射影は \hat{P} と Q とを結ぶ m -測地線により与えられる。
 - 指数型分布族では曲線が η 座標に線形であるとき、 m -測地線となる。
 - m -測地線は \hat{P} で M と直交。直交性はフィッシャー・リーマン情報量で定義。
- 与えられた P に対し $K(Q||P)$ を最小化する $\hat{Q} \in D$ は P の D への e -射影で与えられる。
 - m -射影は P と \hat{Q} とを結ぶ e -測地線により与えられる。 e -測地線は \hat{Q} で D と直交。
 - 指数型分布族では曲線が θ 座標に線形であるとき、 e -測地線となる。
- 観測される部分多様体が D_T^* である場合、 $K(Q^*||P^*) = \sum_{t=1}^T K(Q_t||P_t)$ として各成分 P_t を D_t に e -射影することにより与えられる。各 t に成分的に適用できる。

em(exponential mixture)アルゴリズム

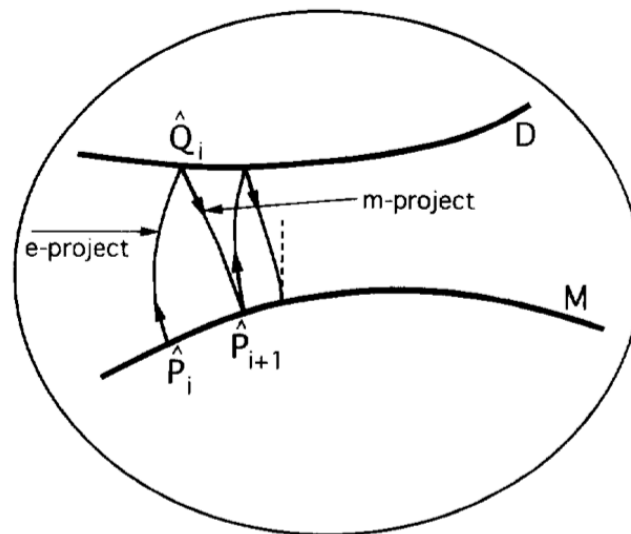
方法

0. 準備. 任意の初期推測値 $\hat{\mathbf{u}}_0$ を選ぶ. 初期推測分布 $\hat{P}_0 \in M$ は $\theta(\hat{\mathbf{u}}_0)$ で与えられる.

1. eステップ. \hat{P}_i を D にe-射影して, $K(Q||\hat{P}_i)$, $Q \in D$ を最小化する \hat{Q}_i を与える.

2. mステップ. \hat{Q}_i を M にm-射影して, $K(\hat{Q}_i||P)$, $P \in M$ を最小化する \hat{P}_{i+1} を与える.

パラメータが収束するまでEステップとMステップを繰り返す.

**FIGURE 9. em algorithm.**

e-射影とm-射影の特性

- ・定理1. \hat{Q} を観測点, M を S 中のモデル多様体とすると, \hat{Q} のm-射影は尤度を最大化する最尤推定量 $\theta(\hat{u}) \in M$ を与える. このm-射影は M がe-平坦であるとき一意である. (Amari,1995)
- ・定理2. D がm-平坦な部分多様体であるとき, P の D へのe-射影は一意である. また D を η 座標で $D = \{\eta \mid \eta = (\eta_v, \eta_h), \eta_v = \check{r}_v, \eta_h \text{は任意}\}$ に分離し, 点 P の θ 座標を $\theta^P = (\theta_v^P, \theta_h^P)$, D へのe-射影を Q^* の η 座標, θ 座標をそれぞれ (η_v^*, η_h^*) , (θ_v^*, θ_h^*) とすると, 以下が成り立つ.

1. Q^* の可視部分 η_v^* は $\eta_v^* = \check{r}_v$ で与えられる.
2. Q^* の θ_h^* はe-射影のもとで普遍に保たれる. $\theta_h^P = \theta_h^*$
3. Q^* における隠れ変数 r_h の条件付き確率は P における確率と等しい.

$$p(r_h | \check{r}_v, P) = p(r_h | \check{r}_v, Q^*)$$

4. Q^* における隠れ変数 r_h の条件付き期待値は P における期待値と等しい.

$$E_P[r_h | \check{r}_v] = E_{Q^*}[r_h | \check{r}_v]$$

e-射影とm-射影の特性

・証明

$Q^* \in D$ なので1の $\eta_v^* = \check{r}_v$ は成り立つ.

KLダイバージェンスは

$$\begin{aligned} K(Q||P) &= E_Q \left[\log \frac{p(\mathbf{r}; \boldsymbol{\theta}_Q)}{p(\mathbf{r}; \boldsymbol{\theta}_P)} \right] \\ &= E_Q [(\boldsymbol{\theta}_v^Q - \boldsymbol{\theta}_v^P) \cdot \mathbf{r}_v + (\boldsymbol{\theta}_h^Q - \boldsymbol{\theta}_h^P) \cdot \mathbf{r}_h] - \psi(\boldsymbol{\theta}_Q) + \psi(\boldsymbol{\theta}_P) \\ &= (\boldsymbol{\theta}_v^Q - \boldsymbol{\theta}_v^P) \cdot \check{r}_v + (\boldsymbol{\theta}_h^Q - \boldsymbol{\theta}_h^P) \cdot \boldsymbol{\eta}_h^Q - \psi(\boldsymbol{\theta}_Q) + \psi(\boldsymbol{\theta}_P) \end{aligned}$$

部分多様体 D において $K(Q||P)$ を最小化するため、微分により

$$\frac{\partial}{\partial \boldsymbol{\eta}_h} K(Q||P) = \frac{\partial \boldsymbol{\theta}_v}{\partial \boldsymbol{\eta}_h} \cdot \check{r}_v + \frac{\partial \boldsymbol{\theta}_h}{\partial \boldsymbol{\eta}_h} \cdot \boldsymbol{\eta}_h + (\boldsymbol{\theta}_h - \boldsymbol{\theta}_h^P) - \frac{\partial \psi(\boldsymbol{\theta}_Q)}{\partial \boldsymbol{\theta}_v} \frac{\partial \boldsymbol{\theta}_v}{\partial \boldsymbol{\eta}_h} - \frac{\partial \psi(\boldsymbol{\theta}_Q)}{\partial \boldsymbol{\theta}_h} \frac{\partial \boldsymbol{\theta}_h}{\partial \boldsymbol{\eta}_h} = 0$$

ここで、 $\frac{\partial \psi(\boldsymbol{\theta}_Q)}{\partial \boldsymbol{\theta}_v} = \boldsymbol{\eta}_v = \check{r}_v$, $\frac{\partial \psi(\boldsymbol{\theta}_Q)}{\partial \boldsymbol{\theta}_h} = \boldsymbol{\eta}_h$ より、2の $\boldsymbol{\theta}_h^P = \boldsymbol{\theta}_h^*$ が成立.

e-射影とm-射影の特性

・ 証明

P における \mathbf{r}_h の条件付き確率は

$$\begin{aligned} p(\mathbf{r}_h | \check{\mathbf{r}}_v, \boldsymbol{\theta}_P) &= \frac{p(\mathbf{r}_h, \check{\mathbf{r}}_v; \boldsymbol{\theta}_P)}{p(\check{\mathbf{r}}_v; \boldsymbol{\theta}_P)} \\ &= \frac{\exp\{\boldsymbol{\theta}_v^P \cdot \check{\mathbf{r}}_v + \boldsymbol{\theta}_h^P \cdot \mathbf{r}_h - \psi\}}{\int \exp\{\boldsymbol{\theta}_v^P \cdot \check{\mathbf{r}}_v + \boldsymbol{\theta}_h^P \cdot \mathbf{r}_h - \psi\} d\boldsymbol{\mu}(\mathbf{r}_h)} \\ &= \exp\{\boldsymbol{\theta}_h^P \cdot \mathbf{r}_h - \tilde{\psi}\} \end{aligned}$$

正規化係数 $\tilde{\psi}$ は $\boldsymbol{\theta}_h^P$ と $\check{\mathbf{r}}_v$ に依存し、 $\boldsymbol{\theta}_v^P$ に依存しない。また、 P と Q で $\boldsymbol{\theta}_h$ は共通であるため、 $p(\mathbf{r}_h | \check{\mathbf{r}}_v, P) = p(\mathbf{r}_h | \check{\mathbf{r}}_v, Q^*)$ が示され、その結果 $E_P[\mathbf{r}_h | \check{\mathbf{r}}_v] = E_{Q^*}[\mathbf{r}_h | \check{\mathbf{r}}_v]$ となる。

EMアルゴリズムの幾何学的構造

・定理1, またはp.21により, Mステップとmステップが同じ手順で $\hat{\mathbf{u}}$ を与えている.

・Eステップの幾何学的な解釈のため, $Q(\boldsymbol{\eta}_v, \boldsymbol{\eta}_h)$ に対し次の変換 $\mathcal{F}: S \rightarrow S$ を定義.

$$\mathcal{F}: (\boldsymbol{\eta}_v, \boldsymbol{\eta}_h) \rightarrow (\boldsymbol{\eta}_v, \mathbf{s}_Q(\boldsymbol{\eta}_v))$$

$$\mathbf{s}_Q(\mathbf{r}_v) = E_Q[\mathbf{r}_h | \mathbf{r}_v]$$

これは, $\boldsymbol{\eta}_v$ は不変, $\boldsymbol{\eta}_h$ は $\mathbf{r}_v = \boldsymbol{\eta}_v$ を条件とした時の \mathbf{r}_h の期待値に置き換える変換.
 $\boldsymbol{\eta}_v = \check{\mathbf{r}}_v$ を固定した時 \mathcal{F} は D を自身に写像する.

・Eステップは P を D にe-射影し, 得られた Q^* を $\mathcal{F}Q^*$ に変換することを意味している.

よってEMアルゴリズムは次のようにかける.

定理3. EMアルゴリズムは次の2つのステップで表される.

1. Eステップ. $K(\mathcal{F}^{-1}(Q) || \hat{P}_i)$, $Q \in D$ を最小化する $\hat{Q}_i \in D$ を与える.
2. Mステップ. $K(\hat{Q}_i || P)$, $P \in M$ を最小化する $\hat{P}_{i+1} \in M$ を与える.

\mathcal{F} が恒等写像の時, EMとemのアルゴリズムが等しい. これは $Q \in D$ での \mathbf{r}_h の条件付き期待値が $Q \in D$ での \mathbf{r}_h の無条件期待値と等しい時である.

EMアルゴリズムの幾何学的構造

定理4. EMとemアルゴリズムは、条件付き期待値 $\mathbf{s}_Q(\mathbf{r}_v) = E_Q[\mathbf{r}_h|\mathbf{r}_v]$ が任意の $Q \in D$ において \mathbf{r}_h に線形であるならば、等価である。

証明

$Q(\boldsymbol{\eta}_v, \boldsymbol{\eta}_h) \in D$ の時、 $\boldsymbol{\eta}_v = E_Q[\mathbf{r}_v]$, $\boldsymbol{\eta}_h = E_Q[\mathbf{r}_h] = E_Q E_Q[\mathbf{r}_h|\mathbf{r}_v] = E_Q[\mathbf{s}_Q(\mathbf{r}_v)]$ とおく。
 $\mathbf{s}_Q(\mathbf{r}_v)$ が一次関数の時、 $\mathbf{s}_Q(\mathbf{r}_v) = \mathbf{a} + \mathbf{B}\mathbf{r}_v$ と定数 \mathbf{a} , 定数行列 \mathbf{B} で記述される。

この時 $\boldsymbol{\eta}_h = E_Q[\mathbf{s}_Q(\mathbf{r}_v)] = \mathbf{s}_Q(E_Q(\mathbf{r}_v)) = \mathbf{s}_Q(\boldsymbol{\eta}_v)$ より等価である。

逆に、 \mathcal{F} が恒等写像であると仮定すると、 $\boldsymbol{\eta}_h = E_Q[\mathbf{s}_Q(\mathbf{r}_v)] = \mathbf{s}_Q(E_Q(\mathbf{r}_v))$ が任意の $Q \in D$ に対して成立。 $Q = \{q(\mathbf{r}_v, \mathbf{r}_h)\}$ を $\boldsymbol{\eta}_v = \boldsymbol{\eta}_{v,0}$ で定義される D_0 に属する任意の分布、 $P = \{p(\mathbf{r}_v, \mathbf{r}_h)\}$ を D_0 への \mathbf{e} -射影が Q に等しくなるような任意の分布とする。 P と Q を結ぶ \mathbf{e} -測地線上の分布の属 Q_t , $0 \leq t \leq 1$ を設定し、 $Q_0 = Q$, $Q_1 = P$ とする。 Q_t は S における指数属で、 Q_t の \mathbf{e} -射影は Q であるため、すべての Q_t は同じ分布条件 $q(\mathbf{r}_v, \mathbf{r}_h)$ を持つ。従って分布 $Q_t = \{q_t(\mathbf{r}_v, \mathbf{r}_h)\}$ は $q_t(\mathbf{r}_v, \mathbf{r}_h) = p(\mathbf{r}_v; t)q(\mathbf{r}_h|\mathbf{r}_v)$ と表される。
 $\mathbf{s}_{Q_t}(\mathbf{r}_v)$ は条件付き分布を通してのみ Q に依存するので、 $\mathbf{s}_{Q_t}(\mathbf{r}_v) = \mathbf{s}_Q(\mathbf{r}_v)$ 。

よって $\int \mathbf{s}_Q(\mathbf{r}_v)p(\mathbf{r}_v; t)d\mathbf{r}_v = \mathbf{s}_Q[\int \mathbf{r}_v p(\mathbf{r}_v, t)d\mathbf{r}_v]$ 。 p は条件 $\int p(\mathbf{r}_v, t)d\mathbf{r}_v = 1$ 下で任意。従って上の式は $\mathbf{s}_Q(\mathbf{r}_v) = \mathbf{a} + \mathbf{B}\mathbf{r}_v$ と同様、 \mathbf{r}_v の一次関数であることを示している。

EMとemアルゴリズムが異なる例：正規分布の乗法モデルの部分観測

- 簡単のために $T = 2$ を考え、写像 \mathcal{F} を探索する.
- D 上の点を $Q = N(\bar{r}, \sigma^2)$ とする.
- $r_v = \frac{1}{2}(x_1 + x_2) = \bar{r}$ が与えられると条件付き確率分布 $p_Q(x_1, x_2 | \bar{r})$ は正規化定数 $c(\bar{r})$ を用いて

$$p_Q(x_1, x_2 | \bar{r}) = c(\bar{r}) \exp \left\{ -\frac{(x_1 - \bar{r})^2 + (x_2 - \bar{r})^2}{2\sigma^2} \right\} \times \delta \left(\frac{x_1 + x_2}{2} - \bar{r} \right)$$

- これを x_2 に関して積分し、 $p_Q(x_1 | \bar{r}) = c \exp \left\{ -\frac{(x_1 - \bar{r})^2}{2\sigma^2} \right\}$. よって \bar{r} がわかると分散は $\frac{\sigma^2}{2}$ に減少する. x_2 も同様.
- よって $r_v = \bar{r}$ を条件とする r_h の条件付き期待値は $s_Q(\bar{r}) = E_Q \left[\frac{x_1^2 + x_2^2}{2} \mid \bar{r} \right] = \bar{r}^2 + \frac{\sigma^2}{2}$ これは無条件の期待値 $\eta_h = \bar{r}^2 + \sigma^2$ とは異なる.
- これより \mathcal{F} は恒等写像ではなく、 $Q = N(\bar{r}, \sigma^2)$ を $\mathcal{F}Q = N\left(\bar{r}, \frac{\sigma^2}{2}\right)$ に写像する.

EMとemアルゴリズムが異なる例：正規分布の乗法モデルの部分観測

- $Q = N(u, u^2) \in M$ の D への e -射影 Q^* は, $\theta_v = \frac{2}{u}$, $\theta_h = \frac{1}{u^2}$ より $T = 2$ で $\eta_v^* = \bar{r}$, $\theta_h^* = \theta_h^P = -u^2$ が成り立つため, $Q^* = N(\bar{r}, u^2)$.
 - 一方, Eステップでは $E_p[r_h | r_v = \bar{r}] = \bar{r}^2 + \frac{u^2}{2}$ より $\hat{Q} = N(\bar{r}, \frac{u^2}{2})$.
- これは, $\mathcal{F}Q^* = \hat{Q}$ を表す.

- emアルゴリズムは $P^* = N(\bar{r}, \bar{r}^2)$ に収束し, 最適な解 $\hat{u}^* = \bar{r}$ を幾何学的に求めるが EMアルゴリズムは最尤推定量 \hat{u} に収束する. $\hat{u} = (\sqrt{3} - 1)\bar{r} \approx 0.73\bar{r}$
これはKLダイバージェンスの最小値ではない

- 一般に, $N > 2$, $Q^* = N(\bar{r}, u^2)$ に対して
$$\mathcal{F}Q^* = N\left(\bar{r}, \frac{T-1}{T} u^2\right)$$

- emの解 $\hat{u}^* = \bar{r}$ に対してEMの解は
$$\hat{u} = \frac{1}{2}(\sqrt{T^2 + 4T} - T)\bar{r} \approx \left(1 - \frac{1}{T}\right)\bar{r}$$

- よって T が大きいとき, 漸近的に等価であり, 現実的には差はない.

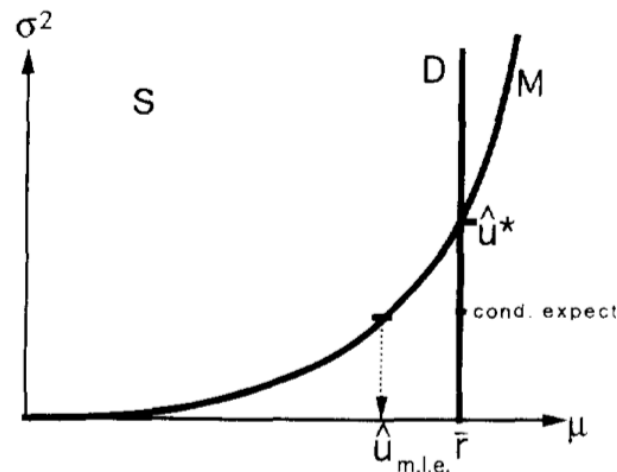


FIGURE 10. Example where EM and em algorithms are different.

EMとemのどちらが自然か

- \mathcal{F} が恒等写像でない場合，異なる推定値を与える.
- EMは最尤推定量を与えるので，統計学的には自然.
- NNで教師データを近似する場合，統計的な推論ではない. その振る舞いはデータ多様体 D にまとめられるため，ダイバージェンスを最小化するemが自然.

- 指数型分布族から関数空間へ枠組みを拡張することにより，二つのアルゴリズムが等価となる.
- また，指数型分布族でも T が大きい場合に漸近的に等価であることを証明する.
- またこの時，条件付き期待値を計算するEステップの方がe-射影より計算量が大きい可能性がある.

EMとemの等価性を保証する線形化トリック

- $s_Q(\mathbf{r}_v)$ が線形である時、等価となる.
- 拡張フレームワークで等価であることを示す.
- y を0か1の値をとる確率変数, $f(y)$ を y に対して線形な関数とする

$$f(y) = \{f(1) - f(0)\}y + f(0)$$

- この方法は有限集合 K 上で値をとる変数であれば, どんな関数でも使える.
 K の要素を m とし, m を指標とする新しいベクトル変数 $k(\mathbf{r})$ を導入する.

$$k(\mathbf{r}) = \{k_m(\mathbf{r}), m \in K, k_m(\mathbf{r}) = \delta(\mathbf{r} - \mathbf{m})\}$$

- $f(\mathbf{r}) = \sum_m f(\mathbf{m})k_m(\mathbf{r})$ より, $f(\mathbf{r})$ は確率変数 $k(\mathbf{r})$ に対して線形

- よって可視確率変数が2値の場合, 可視確率変数が拡張ベクトル $k(\mathbf{r}) = \{k_m(\mathbf{r})\}$ の場合, 条件付き期待値は $k(\mathbf{r})$ に対して線形であり, EMとemが等価.

→2値確率パーセプトロン、エキスパートネットの混合、ボルツマンマシン、正規混合モデルにおいて、EMアルゴリズムとemアルゴリズムが等価であることがわかる

EMとemの等価性を保証する線形化トリック

- ・連続変数の場合にも線形化のトリックを用いる.

$$f(\mathbf{r}) = \int f(\mathbf{m})\delta(\mathbf{r} - \mathbf{m})d\mathbf{m}$$

- ・この場合, 確率変数 \mathbf{r} は一般化された関数に拡張され, \mathbf{r} は一般化関数空間の要素.
- ・ T 個のサンプル $\mathbf{r}_1, \dots, \mathbf{r}_T$ が観測された時, これらは経験分布に要約される.

$$\hat{p}_{emp}(\mathbf{r}) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{r} - \mathbf{r}_t)$$

- ・隠れ変数を関数空間で定式化し, EMとemが関数空間で等価であることを証明する $S = \{p(\mathbf{r}_v, \mathbf{r}_h)\}$ を確率変数 $\mathbf{r} = (\mathbf{r}_v, \mathbf{r}_h)$ の密度関数の関数空間とする. \mathbf{u} で指定されるパラメトリックモデル $M = \{p(\mathbf{r}_v, \mathbf{r}_h; \mathbf{u})\}$ を考え, $\mathbf{r}_{v,T}^* = (\mathbf{r}_{v,1}, \dots, \mathbf{r}_{v,T})$ が観測された時, 関数空間にまとめると

$$\hat{p}_{emp}(\mathbf{r}_v) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{r}_v - \mathbf{r}_{v,t})$$

これは観測データに基づく経験分布.

EMとemの等価性を保証する線形化トリック

$p(\mathbf{r}_v, \mathbf{r}_h) = p(\mathbf{r}_v)p(\mathbf{r}_h|\mathbf{r}_v)$ と分解. $p(\mathbf{r}_v)$ は経験分布で与えられるが $p(\mathbf{r}_h|\mathbf{r}_v)$ は自由.

従って観測データ多様体 D は関数空間において $D = \{\hat{p}_{emp}(\mathbf{r}_v)p(\mathbf{r}_h|\mathbf{r}_v)\}$

$p(\mathbf{r}_v, \mathbf{r}_h; \mathbf{u})$ の D へのe-射影は $K[q(\mathbf{r}_v, \mathbf{r}_h)||p(\mathbf{r}_v, \mathbf{r}_h; \mathbf{u})]$, $q(\mathbf{r}_v, \mathbf{r}_h) \in D$ を最小化するものである. e-射影の Q^* が次式で表せることを示せば良い.

$$\left\{ \begin{array}{l} T = 1 \text{の時} : q^*(\mathbf{r}_v, \mathbf{r}_h) = \hat{p}_{emp}(\mathbf{r}_v)p(\mathbf{r}_h|\mathbf{r}_v; \mathbf{u}) \\ \text{一般的な} T : \prod_{t=1}^T \hat{p}_{emp}(\mathbf{r}_{v,t})p(\mathbf{r}_{h,t}|\mathbf{r}_v; \mathbf{u}) \end{array} \right.$$

隠れた確率変数を関数 $\delta(\mathbf{r}_h - \mathbf{m})$ に拡張すると, その条件付き確率は

$$E[\delta(\mathbf{r}_h - \mathbf{m})|\mathbf{r}_v; \mathbf{u}] = \int \delta(\mathbf{r}_h - \mathbf{m})p(\mathbf{m}|\mathbf{r}_v; \mathbf{u})d\mathbf{m} = p(\mathbf{r}_h|\mathbf{r}_v; \mathbf{u})$$

よって二つのアルゴリズムは同じである.

EMとemの等価性を保証する線形化トリック

定理5. EMとemは関数空間において等価である. その手順は

E(e)ステップ. 条件付き確率 $p(\mathbf{r}_h | \mathbf{r}_v; \hat{\mathbf{u}}_{(i)})$ を求める.

M(m)ステップ. $E_{Q_i}[\log p(\mathbf{r}_h | \mathbf{r}_v; \mathbf{u})] = \sum_{t=1}^T \int p(\mathbf{r}_{h,t} | \mathbf{r}_{v,t}; \hat{\mathbf{u}}_{(i)}) \log p(\mathbf{r}_{h,t} | \mathbf{r}_{v,t}; \mathbf{u}) d\mathbf{r}_{h,t}$ を最大化する \mathbf{u} を探索し, $\hat{\mathbf{u}}_{(i+1)}$ とする.

ただしこの方法では, $\mathbf{r}_{v,1}, \dots, \mathbf{r}_{v,T}$ を統計量にまとめることなく全て保持する必要がある.

例-非線形確率パーセプトロン

• すべてのNNが指数関数的な族だとは限らない.

• f をシグモイド関数 $f(u) = \frac{1}{1+\exp(-u)}$ とする.

• i 番目の隠れユニット, 最終出力はそれぞれ

$$z_i = f(\mathbf{w}_i \cdot \mathbf{x}) + n_i, \quad y = (\mathbf{v} \cdot \mathbf{z}) + n$$

• この時の確率は次のようになる.

$$p(\mathbf{z}|\mathbf{x}; \mathbf{u}) = c \exp \left[-\frac{1}{2\rho^2} \sum \{z_i - f(\mathbf{w}_i \cdot \mathbf{x})\}^2 \right]$$

$$p(y|\mathbf{z}; \mathbf{u}) = c' \exp \left[-\frac{1}{2\rho^2} \{y - f(\mathbf{v} \cdot \mathbf{z})\}^2 \right]$$

従って共同確率分布の対数は

$$\rho^2 l(y, \mathbf{z}|\mathbf{x}; \mathbf{u}) = \sum z_i f(\mathbf{w}_i \cdot \mathbf{x}) + y f(\mathbf{v} \cdot \mathbf{z}) - \frac{1}{2} \{f(\mathbf{v} \cdot \mathbf{z})\} + k(y, \mathbf{z}) - \psi$$

これは $\theta(\mathbf{u}, \mathbf{x})$ と r の形でまとめることができず, $p(y, \mathbf{z}|\mathbf{x}; \mathbf{u})$ は指数型分布族に属さないが, 関数の多様体を導入することで, 情報幾何学とEMアルゴリズムが適用できる.

\mathbf{x} : 入力

\mathbf{z} : 隠れユニットの出力

y : 最終出力

n_i, n : それぞれ $N(0, \sigma^2)$ に従う
独立したランダムノイズ

例-非線形確率パーセプトロン

・確率的出力 z_i, y をそれぞれ期待値に置き換えることで通常のアナログパーセプトロンに還元される。

よって通常が多層パーセプトロンでも確率的手法で学習できる。

・EMアルゴリズムは以下のとおり

Eステップ. 時刻 t の観測値 (y_t, \mathbf{x}_t) に基づき, \mathbf{z}_t の分布を計算

$$\begin{aligned} \mathbf{f}_t &= (f_{1,t}, \dots, f_{k,t}) \\ f_{j,t} &= f(\hat{\mathbf{w}}_j \cdot \mathbf{x}_t) \end{aligned}$$

$$p(\mathbf{z}_t | y_t, \mathbf{x}_t; \hat{\mathbf{u}}_i) = \frac{\exp\left\{-\frac{1}{2\sigma^2} [|\mathbf{z}_t - \mathbf{f}_t|^2 + \{y_t - f(\hat{\mathbf{v}}_i \cdot \mathbf{z}_t)\}^2]\right\}}{\int \exp\left\{-\frac{1}{2\sigma^2} [|\mathbf{z}_t - \mathbf{f}_t|^2 + \{y_t - f(\hat{\mathbf{v}}_i \cdot \mathbf{z}_t)\}^2]\right\} d\mathbf{z}_t}$$

Mステップ. 以下の式を最大化する \mathbf{u} を求め, $\hat{\mathbf{u}}_{i+1}$ とする.

$$\sum_t \int p(\mathbf{z}_t | y_t, \mathbf{x}_t; \hat{\mathbf{u}}_i) \log p(\mathbf{z}_t | y_t, \mathbf{x}_t; \mathbf{u}) d\mathbf{z}_t$$

これは, 最尤推定量に収束する.

- EMとemは時間 t ごとに部分観測 $\mathbf{r}_{v,t}$ または $\mathbf{s}_{v,t}$ が利用できる場合、バッチアルゴリズムではデータ $\mathbf{r}_{v,t}$, $t = 1, \dots, T$ を使用するのに対し、オンライン学習形式で書き直すことができる。
- 学習は一般に収束は遅くなるが、アルゴリズムは単純。
- 環境構造の変化に対してロバスト。

• 全データが S における観測点 $\hat{\boldsymbol{\eta}}$ に集約される場合の学習アルゴリズムを考える。

1.eステップ(Eステップ). 現在の \hat{P}_t の D_{t+1} へのe-射影($\hat{\mathbf{u}}_t$ に基づく $\mathbf{s}_{v,t+1}$ を条件とした $\hat{\mathbf{r}}_{t+1}$ の条件付き期待値)を計算. $\hat{\mathbf{r}}_{t+1}$ が得られ, これを次の式で修正.

$$\hat{\boldsymbol{\eta}}_{t+1} = (1 - \varepsilon_t)\hat{\boldsymbol{\eta}}_t + \varepsilon_t\hat{\mathbf{r}}_{t+1}$$

2.Mステップ. $\hat{\boldsymbol{\eta}}_{t+1}$ から最尤推定量を計算し, $\hat{\mathbf{u}}_{t+1}$ とすると, \hat{P}_{t+1} が求まる.

$\hat{\mathbf{u}}_t$: 時刻 t における推定量
 $\hat{\boldsymbol{\eta}}_t$: 時刻 t における観測点の推定
 $\mathbf{s}_{v,t+1}$: 時刻 $t + 1$ における部分観測点
 ε_t : 定数または $\varepsilon_t \sim \frac{c}{t}$ のような減少数列

- ・最尤推定の計算は通常簡単ではなく、勾配法を用いる。
- ・ $\hat{\boldsymbol{\eta}}_{t+1}$ の対数尤度は $l = \boldsymbol{\theta}(\mathbf{u}) \cdot \hat{\boldsymbol{\eta}}_{t+1} - \psi(\mathbf{u})$ であり、その勾配は

$$\frac{\partial l}{\partial \mathbf{u}} = \mathbf{B}\{\hat{\boldsymbol{\eta}}_{t+1} - \boldsymbol{\eta}(\mathbf{u})\}$$

ただし、 \mathbf{B} は行列 $\mathbf{B} = \frac{\partial \boldsymbol{\theta}}{\partial \mathbf{u}}$ 。

- ・従って増分Mステップにおける勾配法は

$$\hat{\mathbf{u}}_{t+1} = \hat{\mathbf{u}}_t + \varepsilon_t \mathbf{B}\{\hat{\boldsymbol{\eta}}_{t+1} - \boldsymbol{\eta}(\hat{\mathbf{u}}_t)\}$$

よってどの加速度法でも適用可能。

- ・スコアリング法では $\hat{\mathbf{u}}_{t+1} = \hat{\mathbf{u}}_t + \varepsilon_t G_t^* \mathbf{B}\{\hat{\boldsymbol{\eta}}_{t+1} - \boldsymbol{\eta}(\hat{\mathbf{u}}_t)\}$
 G_t^* は $\hat{\mathbf{u}}_t$ におけるフィッシャー情報行列の逆行列

$\hat{\mathbf{u}}_t$: 時刻 t における推定量
 $\hat{\boldsymbol{\eta}}_t$: 時刻 t における観測点の推定
 $\mathbf{s}_{v,t+1}$: 時刻 $t + 1$ における部分観測点
 ε_t : 定数または $\varepsilon_t \sim \frac{c}{t}$ のような減少数列

・観測データが一つの $\hat{\eta}$ にまとめられず、積空間 S_T^* を考慮する必要がある場合

1.eステップ(Eステップ). θ 座標を $(\mathbf{x}_{t+1}, \mathbf{u})$ とするの \hat{P}_t の D_{t+1} へのe-射影($\hat{\mathbf{r}}_{t+1}$ の条件付き期待値)を計算. $\hat{\mathbf{r}}_{t+1}$ が得られ, これを次の式で修正.

2.増分Mステップ. $\hat{\mathbf{u}}_{t+1} = \hat{\mathbf{u}}_t + \varepsilon \mathbf{B}(\mathbf{x}_{t+1}, \hat{\mathbf{u}}_t) \{\hat{\mathbf{r}}_{t+1} - \eta(\hat{\mathbf{u}}_t, \mathbf{x}_t)\}$

この場合, 対数尤度は $l = \sum_t \{\boldsymbol{\theta}(\mathbf{x}_t, \mathbf{u}) \cdot \hat{\mathbf{r}}_t - \psi(\mathbf{x}_t, \mathbf{u})\}$

eステップでは観測変数 $\mathbf{s}_{v,t+1}$ と $\hat{\mathbf{u}}_t$ に基づく \mathbf{x}_{t+1} から $\hat{\mathbf{r}}_{t+1}$ を推定

増分Mステップでは古いデータは破棄し, 最新データ $\hat{\mathbf{r}}_{t+1}$ のみから勾配を計算

$\hat{\mathbf{u}}_t$: 時刻 t における推定量

$\hat{\eta}_t$: 時刻 t における観測点の推定

$\mathbf{s}_{v,t+1}$: 時刻 $t+1$ における部分観測点

ε_t : 定数または $\varepsilon_t \sim \frac{c}{t}$ のような減少数列

例-確率的パーセプトロンの学習

・EMとemは同じ。E(e)ステップにおいて、観測された y_{t+1} を条件とする期待値は

$$\begin{aligned}\hat{r}_{2, \mathbf{k}} &= E[\delta_{\mathbf{k}}(\mathbf{z})] = \text{Prob}\{z = \mathbf{k} | y_{t+1}, \mathbf{u}_t\} \\ &= \frac{\varphi(y_{t+1}, \mathbf{k} \cdot \mathbf{v}_t) \prod \varphi(k_i, \mathbf{x} \cdot \mathbf{w}_{i,t})}{\sum_{\mathbf{k}} \varphi(y_{t+1}, \mathbf{k} \cdot \mathbf{v}_t) \prod \varphi(k_i, \mathbf{x} \cdot \mathbf{w}_{i,t})} \\ \hat{r}_{1, \mathbf{k}} &= \begin{cases} 0, & y_{t+1} = 0 \\ \hat{r}_{2, \mathbf{k}}, & y_{t+1} = 1. \end{cases}\end{aligned}$$

Mステップは最尤推定の計算である。増分アルゴリズムでは

$$l = \boldsymbol{\theta}(\mathbf{u}) \cdot \hat{\mathbf{r}} - \psi(\boldsymbol{\theta}(\mathbf{u}))$$

$\frac{\partial l}{\partial \mathbf{u}}$ を計算し、

$$\begin{aligned}\frac{\partial l}{\partial \mathbf{v}} &= \sum_{\mathbf{k}} \hat{r}_{2, \mathbf{k}} \{y - \varphi(\mathbf{1} \cdot \mathbf{k} \cdot \mathbf{v})\} \mathbf{k}, \\ \frac{\partial l}{\partial \mathbf{w}_i} &= \sum_{\mathbf{k}} \delta(k_i) \hat{r}_{2, \mathbf{k}} \mathbf{x} - \frac{\mathbf{x} \exp(\mathbf{w}_i \cdot \mathbf{x})}{1 + \exp(\mathbf{w}_i \cdot \mathbf{x})}.\end{aligned}$$

$\hat{\mathbf{u}}_{t+1} = \hat{\mathbf{u}}_t + \varepsilon_t \frac{\partial l}{\partial \mathbf{u}}$ で $\hat{\mathbf{u}}_{t+1}$ が更新できる。

例-確率的パーセプトロンの学習

・ 注意点1. ノイズのない (決定論的) ネットワークで $y = 0$ または 1 で例を生成した場合、 $|w_i|$ と $|v|$ は学習によりすべて無限大に発散する。このような場合、正規化によってそれらの大きさを注意深く制御する必要がある。そこで、このような場合にも学習は有効である。

・ 注意点2. $y=0,1$ を出力とする確率的パーセプトロンを用いてアナログ入出力関係 $\bar{y} = f(\mathbf{x})$ を学習する場合、例 $(\bar{y}_1, \mathbf{x}_1), \dots, (\bar{y}_T, \mathbf{x}_T)$ の \bar{y}_t の値はバイナリではなく実数である。このような場合の \bar{y}_t は、 \mathbf{x}_t が入力されたとき、 \bar{y}_t が $y_t=1$ の確率であると解釈される。つまり、 \mathbf{x}_t が多数回入力されたとき、 $y_t=1$ が相対的な頻度 \bar{y}_t で発生すると解釈する。条件付き期待値は次式で与えられる。

$$\begin{aligned}\hat{r}_{1,k} &= \bar{y}_{t+1} E[\delta_{\mathbf{k}}(\mathbf{z}) | y_{t+1} = 1, \mathbf{u}_t], \\ \hat{r}_{2,k} &= \bar{y}_{t+1} E[\delta_{\mathbf{k}}(\mathbf{z}) | y_{t+1} = 1, \mathbf{u}_t] \\ &\quad + (1 - \bar{y}_{t+1}) E[\delta_{\mathbf{k}}(\mathbf{z}) | y_{t+1} = 0, \mathbf{u}_t].\end{aligned}$$

- **EM**アルゴリズムと**em**アルゴリズムの情報幾何学的知見が構築された.
- 両者は**T**が大きくなると漸近的に等価になり, 実用的なほとんどのケースで等価.
- 拡張関数空間においても等価
- これにより**EM**の学習過程を定式化することが可能になった.