

階層型ネットワークの階層型ネットワーク

浅川伸一 <asakawa@twcu.ac.jp>

1 はじめに

Mixture of experts (以下 ME と略記) と呼ばれる階層型のネットワークを紹介します。ME は入力データ空間をいくつかの小領域に分割し、その分割された各領域に対してひとつのニューラルネットワークを割り当てることによって複雑な問題の解を求めるため手法です。複雑な問題を分割して小領域に区切ることによって、1つの大きなニューラルネットワークを学習させるよりも効率の良い学習をさせようというモデルです。このような手法を分割統治 divide-and-conquer と言ったりします。ME における学習とは、入力空間の分割の仕方と、各小領域に属するデータに対する最適な答えを見つけ出すことです。

分割統治 は科学における一般原理であるといつて良いでしょう。この分割統治を自動的に行ない学習させようと言うのが ME の発想です。ME はローカルネットワークとゲートネットワークから成り立っています。ME は教師あり学習の一手法です。

2 エキスパートの階層

入力ベクトルが m 次元の実数値をとるベクトル \mathcal{R}^m で、出力が n 次元の実数値をとるベクトル \mathcal{R}^n で定義された要素であるとしします。観測されたデータを $\mathcal{X} = \{x^{(t)}, y^{(t)}\}$ としします。ME の木構造において非末端部にはゲーティングネットワークがあり、1つの数値(確率と考えると良い)を出力します。一方、末端部分にはエキスパートネットワークがあります。各エキスパートネットワークは入力ベクトル x に対して出力ベクトル μ_{ij} を出力します。各々の出力ベクトルは木を登って行きゲーティングネットワークによって混合されます。

エキスパートネットワーク (i, j) は入力ベクトル x に対する出力としてベクトル μ_{ij} を出力します。

$$\mu_{ij} = f(U_{ij}x) \quad (1)$$

ここで U_{ij} は結合係数行列で、入力 x は切片項を表す常に 1 を出力する固定の要素を含みます。関数 f のことを入力と出力を結びつけるという意味で

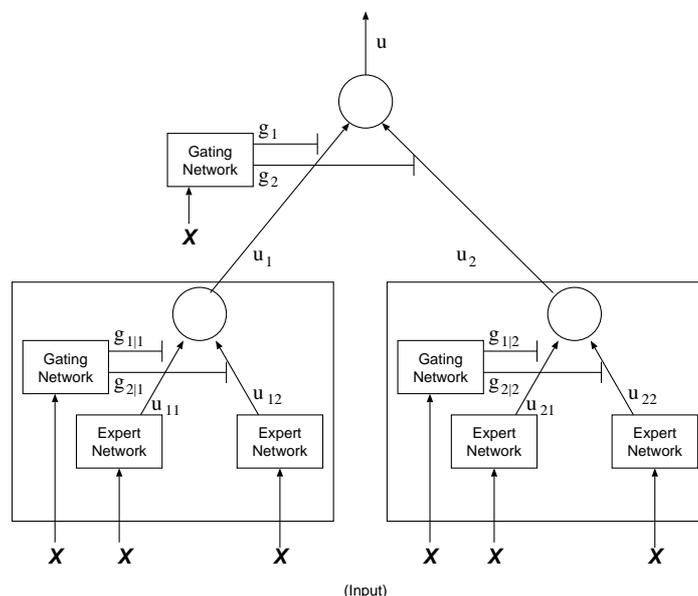


図 1: 2 段階の ME

リンク関数と言います。 f は問題によって線形でもシグモイド関数でも、あるいはガウシアン関数でもなんでも良いのです。

ゲーティングネットワークは線形であり ξ を変数とする以下の式のように表されます。

$$\xi_i = \mathbf{v}_i^T \mathbf{x} \quad (2)$$

ここで \mathbf{v}_i は結合係数ベクトルです。このとき最上位のゲーティングネットワークの i 番目の出力は ξ_i のソフトマックス関数として定義されます。

$$g_i = \frac{e^{\xi_i}}{\sum_k e^{\xi_k}} \quad (3)$$

g_i は正であり、すべての g_i について合計すると必ず 1 になります。すなわちソフトマックス関数は入力空間を「柔らかく」分割するものだといえます (第 3 節参照)。

同じようにして下位のゲーティングネットワークも線形であり、出力 ξ_{ij} は以下のように定義されます。

$$\xi_{ij} = \mathbf{v}_{ij}^T \mathbf{x} \quad (4)$$

このとき

$$g_{j|i} = \frac{e^{\xi_{ij}}}{\sum_k e^{\xi_{ik}}} \quad (5)$$

は 2 段目の i 番目のゲーティングネットワークにおける j 番目のゲーティングネットワークの出力です。ここでも $g_{j|i}$ は正であり各 \mathbf{x} に対して足し合わ

せると 1 になります。非末端部 (ノード) における出力は、そのノードの下にあるエキスパートネットワークの出力に各ゲーティングネットワークの重みをかけた出力になります。すなわち第 2 層のノード部における出力は

$$\mu_i = \sum_j g_{j|i} \mu_{ij} \quad (6)$$

となり、木の最上位層ゲーティングネットワークの出力は

$$\mu = \sum_i g_i \mu_i = \sum_i g_i \sum_j g_{j|i} \mu_{ij} \quad (7)$$

となります。\$g\$ も \$\mu\$ も入力 \$x\$ に依存しています、故に、全体の出力は入力ベクトル \$x\$ の非線型関数となります。\$g_i\$ は \$v_i\$ と \$x\$ とに依存する関数であることを明記するために \$g_i(v_i, x)\$ と表記することもあります。

3 リッジ回帰

ソフトに分割された領域とは、領域間に重なりがあるという意味です。この重複の性質を理解するために 2 つのエキスパートからなる 1 つの階層を考えてみることにします。この場合、ゲーティングネットワークは \$g_1\$ と \$g_2\$ との 2 つの出力を持ちます。ゲーティング出力 \$g_1\$ は

$$g_1 = \frac{e^{\xi_1}}{e^{\xi_1} + e^{\xi_2}} \quad (8)$$

$$= \frac{1}{1 + e^{-(v_1 - v_2)x}} \quad (9)$$

と表されます。これは \$v_1 - v_2\$ というベクトルの方向によって決まるシグモイド関数です。ゲーティング出力 \$g_2\$ は \$1 - g_1\$ に等しいこととなります。あ

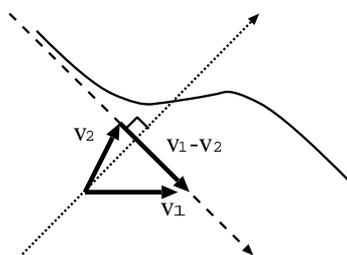


図 2: ソフトマックスとリッジ

る \$x\$ が与えられると全体の出力 \$\mu\$ は \$g_1 \mu_1 + g_2 \mu_2\$ によって与えられます。

これはエキスパートの重みづけ平均となり、重みはリッジ関数の値によって決まります。 $g_1 = g_2 = 1/2$ というリッジに沿って 2 つのエキスパートは等しく貢献します。このリッジから離れるに従って 1 つのエキスパートの影響が大きくなり、他方の影響は小さくなります。

$$\begin{aligned} g_1 = g_2 = \frac{1}{2} &= \frac{1}{1 + e^{-(v_1 - v_2)x}} \\ 1 &= \frac{2}{1 + e^{-(v_1 - v_2)x}} \\ e^{-(v_1 - v_2)x} &= 1 \\ -(v_1 - v_2)x &= 0 \end{aligned}$$

という計算により $v_1 - v_2$ というベクトルに直交するベクトルの方向では 2 つのエキスパートの影響は等しくなります。リッジを横切る滑らかさの度合は $v_1 - v_2$ というベクトルの大きさによって決まります。もし $v_1 - v_2$ が大きければリッジ関数は 2 つの領域をハッキリと分離することになり、エキスパートの重み付された出力は線形に近くなります。もし、 $v_1 - v_2$ が小さいのなら各エキスパートはリッジの両側において意味のある程度まで貢献し、なめらかなマップになります。 $v_1 - v_2$ の値が大きい程分離の程度が大きくなり、この重みが 0 になる極限では両エキスパートは等しく貢献することになります。

4 ME の確率的解釈

ME の出力ベクトル y を、入力データベクトル x とデータの密度関数を表すパラメータ θ とを用いて条件付確率として定式化することができます。

$$P(y|x, \theta) = \sum_i g_i(v_i, x) \sum_j g_{j|i}(v_{ij}, x) P(y|x, \theta_{ij}) \quad (10)$$

ここで θ_{ij} は分布の密度関数を定めるパラメータベクトルです。例えば $P(y|x, \theta)$ が正規分布に従うという仮定の下では θ の成分には平均と分散が含まれます。データが多次元正規分布に従い、分散共分散行列が $\Sigma = \sigma^2 I$ (ただし I は単位行列) となるならば、データ x と θ とが与えられたときの密度関数は次のように書くことができます。

$$p(y|x, \theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \sum_i g_i \sum_j g_{j|i} e^{-(1/2\sigma^2)(y - \mu_{ij})^T (y - \mu_{ij})} \quad (11)$$

分散パラメータ σ^2 は μ_{ij} を中心とする円の大きさ (半径) を定めるものと解釈することができます。 $\sigma^2 \rightarrow 0$ の極限ではディラックのデルタ関数となります。

y が離散的な 2 値関数、すなわち多次元ベルヌーイ試行ならば

$$p(y|x, \theta) = \sum_i g_i \sum_j g_{j|i} u_{ij}^y (1 - u_{ij})^{1-y} \quad (12)$$

となります。

5 パラメータの推定

5.1 事後確率

ME ニューラルネットワークの学習を考えるために ME の木構造のノードにおける事後確率 h を定義します。 g_i と $g_{j|i}$ とは事前確率と呼ぶことにします。なぜなら $g_i, g_{j|i}$ は入力データベクトル \mathbf{x} だけで決まり、対応する出力ベクトル \mathbf{y} の知識なしに定まるからです。事後確率は入力ベクトル \mathbf{x} と出力ベクトル \mathbf{y} の両方が分かったとき定義されます。ベイズの定理を用いて i 番目のノードの事後確率は次のように定義されます。

$$h_i = \frac{g_i \sum_j g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})}{\sum_i g_i \sum_j g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})} \quad (13)$$

$$h_{j|i} = \frac{g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})}{\sum_j g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})} \quad (14)$$

g_i , と $g_{j|i}$ は (3), (5) 式で定義されています。

また、 h_i と $h_{j|i}$ との積、すなわち結合事後確率 h_{ij} を定義しておきます。この量はエキスパートネットワーク (i, j) がデータを生成したと見なすことができる確率です。

$$h_{ij} = \frac{g_i g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})}{\sum_i g_i \sum_j g_{j|i} P_{ij}(\mathbf{y}|\mathbf{x})} \quad (15)$$

5.2 尤度

N 個のデータ集合 $\mathcal{X} = \{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}_1^N$ から次の対数尤度関数が得られます。

$$l(\theta; \mathcal{X}) = \sum_t \log \sum_i g_i^{(t)} \sum_j g_{j|i}^{(t)} P_{ij}(\mathbf{y}^{(t)}|\mathbf{x}^{(t)}) \quad (16)$$

密度関数 P がガウシアンで共分散行列が単位行列であり、リンク関数が恒等関数である場合を考えます。このとき $l(\theta; \mathcal{X})$ を各々のパラメータで微分することで以下のような結合係数行列の勾配上昇学習則を得ます。

$$\Delta U_{ij} = \rho \sum_t h_i^{(t)} h_{j|i}^{(t)} (\mathbf{y}^{(t)} - \boldsymbol{\mu}_{ij}^{(t)}) \mathbf{x}^{(t)T} \quad (17)$$

ここで ρ は学習率です。一番上のゲーティングネットワークにおける i 番目の結合係数ベクトルの勾配上昇学習則は

$$\Delta \mathbf{v}_i = \rho \sum_t (h_i^{(t)} - g_i^{(t)}) \mathbf{x}^{(t)} \quad (18)$$

で与えられます。 i 番目の下位レベルのゲーティングネットワークにおける j 番目の結合係数ベクトルの勾配上昇学習則は

$$\Delta v_{ij} = \rho \sum_t h_i^{(t)} \left(h_{j|i}^{(t)} - g_{j|i}^{(t)} \right) \mathbf{x}^{(t)} \quad (19)$$

式 (17),(18),(19) で与えられるアルゴリズムはバッチ学習アルゴリズムです。対応するオンラインアルゴリズムは単純に加算記号を落して、刺激提示毎にパラメータを更新することで得られます。例えば

$$U_{ij}^{(t+1)} = U_{ij}^{(t)} + \rho h_i^{(t)} h_{j|i}^{(t)} \left(\mathbf{y}^{(t)} - \boldsymbol{\mu}^{(t)} \right) \mathbf{x}^{(t)T} \quad (20)$$

は t 番目の刺激パターンに基づく (i, j) 番目のエキスパートネットワークの結合係数に対する確率的更新則です。

5.3 EM アルゴリズム

EM アルゴリズムによって ME ネットワークのパラメータを求めることを考えます。まず尤度関数が簡単になるような隠れ変数を定義します。変数 z_i はすべての i の中のたった 1 つの z_i だけが 1 となり、他の全ては 0 であるとします。同様に変数 $z_{j|i}$ はたった 1 つの $z_{j|i}$ だけが 1 であり、他はすべて 0 であるとします。これらの隠れ変数は確率モデルにおける決定に対応したラベルであると解釈できます。また z_i と $z_{j|i}$ との積である z_{ij} を定義しておきます。 z_{ij} は特定のエキスパートネットワークを示すラベルであると解釈できます。もし、ラベル $z_i, z_{j|i}, z_{ij}$ が分かれば尤度最大化問題は各エキスパートネットワークの回帰問題に分離され、ゲーティングネットワークについては分割された集合の分類問題に切り離され、これらの問題は互いに独立に解くことができます。もちろん隠れ変数は未知なのですが、隠れ変数と観測可能なデータとを結びつける確率モデルを特定することができます。この確率モデルは z_{ij} の項を用いて

$$P \left(\mathbf{y}^{(t)}, z_{ij}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta} \right) = g_i^{(t)} g_{j|i}^{(t)} P_{ij}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}) \quad (21)$$

$$= \prod_i \prod_j \left\{ g_i^{(t)} g_{j|i}^{(t)} P_{ij}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}) \right\}^{z_{ij}^{(t)}} \quad (22)$$

と書けます。この確率の対数をとると次の完全データ尤度が求められます。

$$l_c(\boldsymbol{\theta}; \mathcal{Y}) = \sum_t \sum_i \sum_j z_{ij}^{(t)} \log \left\{ g_i^{(t)} g_{j|i}^{(t)} P_{ij}(\mathbf{y}^{(t)}) \right\} \quad (23)$$

$$= \sum_t \sum_i \sum_j z_{ij} \left\{ \log g_i^{(t)} + \log g_{j|i}^{(t)} + \log P_{ij}(\mathbf{y}^{(t)}) \right\} \quad (24)$$

上式で表現された完全データ尤度と (16) で定義された不完全データ尤度の関係に注意して下さい。指示変数 z_{ij} を用いることで対数を加算記号の内側

に持って来ることが可能になりました。このことによって最大化問題を簡単
にすることができます。完全データ尤度の期待値をとることで EM アルゴリ
ズムの E ステップを

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)}) = \sum_t \sum_i \sum_j h_{ij}^{(t)} \left\{ \log g_i^{(t)} + \log g_{j|i}^{(t)} + \log P_{ij}(\mathbf{y}^{(t)}) \right\} \quad (25)$$

と定義します。ここで以下の事実を使いました。

$$E \left[z_{ij}^{(t)} | \mathcal{X} \right] = P \left(z_{ij}^{(t)} = 1 | \mathbf{y}^{(t)}, \mathbf{x}^{(t)}, \boldsymbol{\theta}^{(p)} \right) \quad (26)$$

$$= \frac{P \left(\mathbf{y}^{(t)} | z_{ij} = 1, \mathbf{x}^{(t)}, \boldsymbol{\theta}^{(p)} \right) P \left(z_{ij}^{(t)} = 1 | \mathbf{x}^{(t)}, \boldsymbol{\theta}^{(p)} \right)}{P \left(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta}^{(p)} \right)} \quad (27)$$

$$= \frac{P \left(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta}_{ij}^{(p)} \right) g_i^{(t)} g_{j|i}^{(t)}}{\sum_i g_i^{(t)} \sum_j g_{j|i}^{(t)} P \left(\mathbf{y}^{(t)} | \mathbf{x}^{(t)}, \boldsymbol{\theta}^{(p)} \right)} \quad (28)$$

$$= h_{ij}^{(t)} \quad (29)$$

すなわち z_{ij} の期待値は事後確率 h_{ij} になります。同様にして $E \left[z_i^{(t)} | \mathcal{X} \right] = h_i^{(t)}$ 、 $E \left[z_{j|i}^{(t)} | \mathcal{X} \right] = h_{j|i}^{(t)}$ が導出できます。

EM アルゴリズムの M ステップでは $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{(p)})$ をエキスパートネットワー
クパラメータとゲーティングネットワークパラメータに関して最大化する必
要があります。式 (25) を見ると、エキスパートネットワークパラメータは
 $h_{ij}^{(t)} \log P(\mathbf{y}^{(t)})$ をとおしてだけ Q に影響し、ゲーティングネットワークパラ
メータは $h_{ij}^{(t)} \log g_i^{(t)}$ と $h_{ij}^{(t)} \log g_{j|i}^{(t)}$ とをとおしてだけ Q に影響するこ
とが分かります。従って M ステップは以下のような個別の最大化問題に帰結し
ます。

$$\boldsymbol{\theta}_{ij}^{(p+1)} = \arg \max_{\boldsymbol{\theta}_{ij}} \sum_t h_{ij}^{(t)} \log P_{ij}(\mathbf{y}^{(t)}) \quad (30)$$

$$\mathbf{v}_i^{(p+1)} = \arg \max_{\mathbf{v}_i} \sum_t \sum_k h_k^{(t)} \log g_k^{(t)} \quad (31)$$

$$\mathbf{v}_{ij}^{(p+1)} = \arg \max_{\mathbf{v}_{ij}} \sum_t \sum_k h_k^{(t)} \sum_l h_{l|k}^{(t)} \log g_{l|k}^{(t)} \quad (32)$$

個々の最大化問題はそれ自身尤度最大化問題です。式 (30) は単に密度関数 P_{ij}
の重みつき尤度最大化問題です。