

HMM の学習アルゴリズムの並列化に関する一考察

広島大学大学院工学研究科 河合 理恵, 岡村 寛之, 土肥 正
Rie Kawai, Hiroyuki Okamura and Tadashi Dohi
Graduate School of Engineering
Hiroshima University, Japan

1 はじめに

HMM (Hidden Markov Model; 隠れマルコフモデル) は時系列データを解析する統計モデルの一つであり, パラメータ未知のマルコフ過程を仮定している. 確率によって動作する非決定有限オートマトンであり, 観測データから内部状態の推移を一意に特定できない. そのために隠れマルコフモデルと呼ばれている. HMM では観測された時系列データからモデル内部のパラメータを推定する. 連続的かつ伸縮しうる信号列のパターン抽出などに適しており, 音声認識をはじめとして DNA 配列設計, 画像認識にも利用されている [1, 2, 3].

HMM の学習は, 観測データから内部のパラメータを推定する作業に対応づけられる. 代表的な学習アルゴリズムに, Baum らが提案した Baum-Welch アルゴリズム (BW) がある [4]. BW は EM 法 (Expectation-Maximization method) を用いてパラメータの最尤推定値を求める学習アルゴリズムであるが, 計算量が観測データの大きさに比例して増加する. そのため, 長い系列データを扱う場合などに膨大な学習時間を要する場合がある.

一般に数値計算の計算時間短縮は, アルゴリズムの改良による理論的な時間計算量の減少と, コーディングによる工夫やハードウェア化によって処理効率を上げる手法 [5] がある. コーディングによる計算時間の短縮は, BLAS (Basic Linear Algebra Subprograms) などのようなアーキテクチャに最適化されたライブラリを利用することで簡単に計算速度を向上することができるため, 本質的なアルゴリズムの改良なしで手軽に利用できる.

一方, 近年の CPU におけるコア数の増加と OpenMP などのメモリ共有型の並列化ライブラリの整備は「低コストの並列化」と言うこれまでとは異なる計算時間短縮の潮流を生み出している. メモリ共有型の並列化は手軽に利用できるが, 既存のアルゴリズムを単純に並列化するだけでは十分な効果を得ることができないことが多い. そこで本稿では, HMM の学習アルゴリズムに対して, 一般化 EM 法 (Generalized EM; GEM) [6] と呼ばれる手法を適用する.

本稿では, 第 2 章は HMM の定義について説明し, 第 3 章では EM 法から導出される Baum-Welch アルゴリズム (BW) を紹介する. 第 4 章では GEM の原理と, GEM を用いた HMM の学習アルゴリズムについて説明する. 第 5 章では, 2 つのアルゴリズムを計算時間と学習の効率性の観点から比較する. 最後に第 6 章では, 今回得られた数値例から得た結論と今後の課題について述べる.

2 HMM の記述

本稿では離散型 HMM について考える. 内部状態数を M , シンボル数を L とする. 内部状態の集合を $S = \{s_1, \dots, s_M\}$, 出力シンボルの集合を $v = \{v_1, \dots, v_L\}$ と表す. 次に, HMM に対して, 内部状態の推移確率行列 $A = [a_{i,j}]$ を定義する. また, シンボル $v_l, l = 1, \dots, L$ の出力に関する確率行列を $B(v_l)$ とする. シンボルの出力は内部状態に依存するため, 対角要素に $b_i(v_l)$ を持つ対角行列として $B(v_l)$ を仮定する. $b_i(v_l)$ は内部状態 s_i でシンボル v_l を出力する確率を表している. また HMM の初期状態確率を表現する確率ベクトルを $\pi = \{\pi_1, \dots, \pi_M\}$, 吸収状態を表すベクトルを $\xi = \{\xi_1, \dots, \xi_M\}^T$ とする. このとき, π_i は内部状態 s_i を初期状態として選択する確率, ξ_i は内部状態 s_i が吸収状態であれば 1, 吸収状態でない場合は 0 である.

3 Baum-Welch アルゴリズム

3.1 EM 法の概要

EM 法は期待値最大化法とも呼ばれ, 観測データからモデルパラメータの最尤推定値を算出するのに用いられる. 一般的に観測データを D , 未観測データ (潜在変数) を Z とする. このとき, 観測データに対する

対数尤度をモデルパラメータベクトル θ を使って表現すると,

$$\log P(D; \theta) = \log \sum_{\mathbf{Z}} P(D, \mathbf{Z}; \theta) \quad (1)$$

となる. ここで $P(\cdot)$ は任意の密度または確率関数である. 原理的に最尤推定値はこの式を最大化する θ を求めればよい. しかしながら, 式(1)を最大にする θ を陽に得ることは困難である. そこでEM法では完全データ (D, \mathbf{Z}) に対する対数尤度関数の期待値を最大にする θ を求めることを考える. 具体的に未観測データに対する条件付分布 $P(\mathbf{Z}|D; \theta)$ が得られたとすると, 完全データに対する期待対数尤度は,

$$Q(\theta|\theta') = \sum_{\mathbf{Z}} P(\mathbf{Z}|D; \theta') \log P(D, \mathbf{Z}; \theta) \quad (2)$$

となる. ここで θ' は暫定的に与えられるパラメータベクトルである. 式(2)で表される完全データに対する期待対数尤度は Q 関数と呼ばれる. EM法では暫定パラメータ θ' が与えられたもとの Q 関数を算出するステップ (Eステップ) と, それを最大にする θ を求め, 新たなパラメータとして更新するステップ (Mステップ) を繰り返すことで最尤推定値を算出する.

3.2 EM法を用いたHMMパラメータの推定 (Baum-Welchアルゴリズム)

観測データとして, シンボル系列 $D = \{x_1, \dots, x_K\}$ が得られた時のHMMパラメータ推定を考える. 未観測データ (潜在変数) として $\mathbf{Z} = \{z_1, \dots, z_{K+1}\}$, $z_k = (z_{k,1}, \dots, z_{k,M})$ を導入する. $z_{k,i}$ は時刻 k においてシンボル x_k を出力する直前の内部状態が i であるかどうかをあらわす指標であり, 内部状態が i ならば1, それ以外は0である. これを用いると完全データ (D, \mathbf{Z}) に対する対数尤度は,

$$\log P(D, \mathbf{Z}) = \sum_{i=1}^M z_{1,i} \log \pi_i + \sum_{k=1}^K \sum_{i=1}^M \sum_{j=1}^M z_{k,i} z_{k+1,j} \log (I(x_k = v_i) b_i(v_i) a_{i,j}) + \sum_{i=1}^M z_{K+1,i} \log \xi_i \quad (3)$$

となる. $P(\mathbf{Z}|D; \theta')$ による期待演算子を $E[\cdot|D; \theta']$ とすると, 式(2)および(3)より,

$$Q(\theta|\theta') = \sum_{i=1}^M E[z_{1,i}|D; \theta'] \log \pi_i + \sum_{k=1}^K \sum_{i=1}^M \sum_{j=1}^M E[z_{k,i} z_{k+1,j}|D; \theta'] \log (I(x_k = v_i) b_i(v_i) a_{i,j}) + \sum_{i=1}^M E[z_{K+1,i}|D; \theta'] \log \xi_i \quad (4)$$

を得る. ここで $\theta = (\pi, \mathbf{A}, \mathbf{B}, \xi)$ である. このとき式(4)を最大化するパラメータはラグランジュ未定乗数法から,

$$a_{i,j} = \frac{\sum_{k=1}^K E[z_{k,i} z_{k+1,j}|D; \theta']}{\sum_{k=1}^K \sum_{j=1}^M E[z_{k,i} z_{k+1,j}|D; \theta']}, \quad b_i(v_l) = \frac{\sum_{k=1}^K I(z_k = v_l) E[z_{k,i}|D; \theta']}{\sum_{k=1}^K E[z_{k,i}|D; \theta']}, \quad \pi_i = E[z_{1,i}|D; \theta'] \quad (5)$$

と陽に得られる (Mステップにおける更新式). 一方, $E[z_{k,i} z_{k+1,j}|D; \theta']$ は,

$$\begin{aligned} E[z_{k,i} z_{k+1,j}|D; \theta] &= \sum_{\mathbf{Z}} z_{k,i} z_{k+1,j} P(\mathbf{Z}|D; \theta) \\ &= \frac{1}{P(D)} [\pi \mathbf{B}(x_1) \mathbf{A} \mathbf{B}(x_2) \mathbf{A} \cdots \mathbf{B}(x_{k-1}) \mathbf{A}]_i b_i(x_k) a_{i,j} [\mathbf{B}(x_{k+1}) \mathbf{A} \cdots \mathbf{B}(x_K) \mathbf{A} \xi]_j \end{aligned} \quad (6)$$

によって求めることができる. ここで α_k, β_k を

$$\alpha_k = \pi \mathbf{B}(x_k) \mathbf{A} \mathbf{B}(x_{k+1}) \mathbf{A} \cdots \mathbf{B}(x_K) \mathbf{A}, \quad \beta_k = \mathbf{B}(x_k) \mathbf{A} \mathbf{B}(x_{k+1}) \mathbf{A} \cdots \mathbf{B}(x_K) \mathbf{A} \xi \quad (7)$$

とすると, $E[z_{k,i}|D]$ と $E[z_{1,i}|D]$ は,

$$E[z_{k,i}|D] = \frac{1}{P(D)} [\alpha_{k-1}]_i [B(x_k) A \beta_{k+1}]_i, \quad (8)$$

$$E[z_{1,i}|D] = \frac{\pi_i [\beta_1]_i}{P(D)} \quad (9)$$

によって得ることができる. ここで新たな変数とその期待値を導入する,

$$E[N_{i,j}|D] = \sum_{k=1}^K E[z_{k,i} z_{k+1,j}|D], \quad E[Y_{i,l}|D] = \sum_{k=1}^K I(x_k = v_l) E[z_{k,i}|D], \quad E[R_i|D] = E[z_{1,i}|D]. \quad (10)$$

$E[N_{i,j}|D]$ は時刻 0 から時刻 K までの間に, 内部状態 s_i が s_j へ推移した回数の期待値であり, 分子を $EN_{i,j}$ で表す. $E[Y_{i,l}|D]$ は時刻 0 から時刻 K までの間に, 内部状態 s_i でシンボル v_l が出力された回数の期待値であり, 分子を $EY_{i,l}$ で表す. $E[R_i|D]$ は HMM が内部状態 s_i で開始されるかどうかの期待値であり分子を ER_i で表す. 実際の BW によるパラメータ推定では, これらの確率変数を用いて以下の 4 つのステップを繰り返すことで推定値を算出する.

Step1 時刻 0 および各時刻 $k = 1, 2, \dots, K$ に対して前向き尤度 α_k を求める.

$$\alpha_0 = \pi, \quad \alpha_k = \alpha_{k-1} B(x_k) A. \quad (11)$$

Step2 時刻 $K + 1$ および各時刻 $k = K, K - 1, \dots, 1$ に対して後ろ向き尤度 β_k を求める.

$$\beta_{K+1} = \xi, \quad \beta_k = B(x_k) A \beta_{k+1}. \quad (12)$$

Step3 確率変数の期待値の分子を求める.

$$EN_{i,j} = \sum_{k=1}^K [\alpha_{k-1}]_i b_i(x_k) a_{i,j} [\beta_{k+1}]_j, \quad (13)$$

$$EY_{i,l} = \sum_{k=1}^K I(x_k = v_l) [\alpha_{k-1}]_i [B(v_l) A \beta_{k+1}]_i, \quad (14)$$

$$ER_i = \pi_i [\beta_1]_i \quad (15)$$

Step4 パラメータを更新する.

$$a_{i,j} = \frac{EN_{i,j}}{\sum_{j=1}^M EN_{i,j}}, \quad b_i(l) = \frac{EY_{i,l}}{\sum_{l=1}^L EY_{i,l}}, \quad \pi_i = \frac{ER_i}{\sum_{i=1}^M ER_i}. \quad (16)$$

ここでは, 各期待値の分母は Step4 における除算で約分されるので, 算出していない. 実際には暫定的なパラメータから始めて, Step1 から Step4 をパラメータまたは尤度が収束するまで繰り返す.

最後に BW の並列化について考える. 並列化を行うためにはデータ依存の除去が重要な課題である. データ依存とはループ内での計算において, 一つ前の反復の値が次の反復に影響を及ぼすような場合を指す. 例えば Step1 の計算において α_{k-1} を求めなければ, α_k を求めることができない. この関係を除去することが並列計算には必要不可欠である.

図 1 は BW の並列化の概念図である. BW は前向き確率と後ろ向き確率の計算にデータ依存を含んでいるため Step1 と Step2 は逐次実行する必要がある. この部分では BW は高々 2 並列での計算が限界である. Step3 は足し算の繰り返しとなっており, データ依存が存在しない. そのためスレッド数に応じた並列計算が可能となる. ステップ 4 ではステップ 3 の結果を用いて除算を行ってパラメータを更新する.

4 一般化 EM 法による並列学習アルゴリズム

4.1 一般化 EM 法の原理

一般化 EM 法 (Generalized EM, GEM) は EM 法における事後分布に対して近似を適用するものである. パラメータ ψ をもつ新たな分布 $\tilde{P}(Z|\psi)$ を導入し $P(Z|D; \theta)$ の近似として用いる. このとき, 観測データに

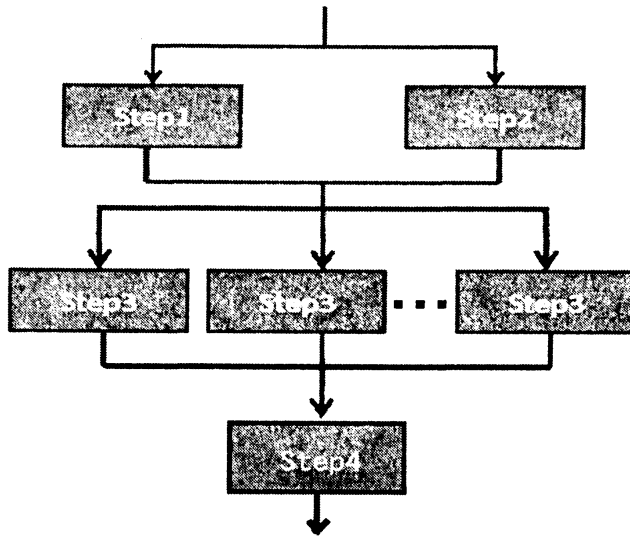


図 1: BW の並列化の概念図.

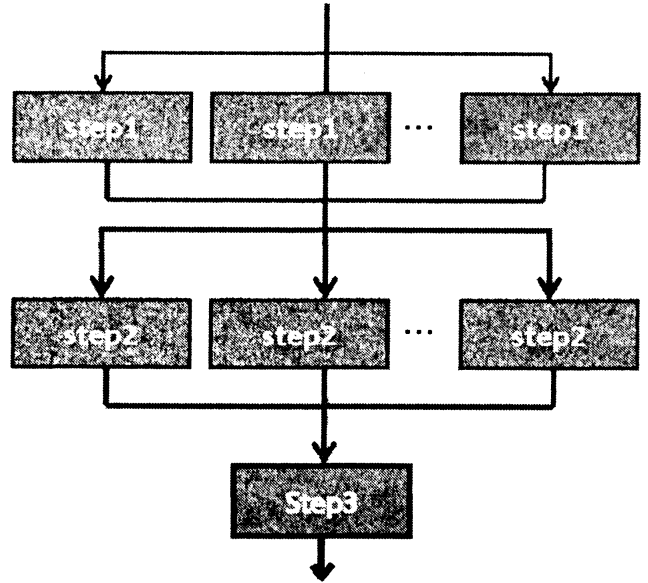


図 2: GEM の並列化の概念図.

対する対数尤度の下限值として以下の結果を得る.

$$\begin{aligned} \log P(\mathbf{D}; \theta) &= \log \sum_{\mathbf{Z}} \tilde{P}(\mathbf{Z}|\psi) \frac{P(\mathbf{D}, \mathbf{Z}; \theta)}{\tilde{P}(\mathbf{Z}|\psi)} \\ &\geq \sum_{\mathbf{Z}} \tilde{P}(\mathbf{Z}|\psi) \log \frac{P(\mathbf{D}, \mathbf{Z}; \theta)}{\tilde{P}(\mathbf{Z}|\psi)} \equiv F(\tilde{P}, \theta). \end{aligned} \quad (17)$$

単純に下限値 $F(\tilde{P}, \theta)$ と対数尤度の差を考えると,

$$\log P(\mathbf{D}, \theta) - F(\tilde{P}, \theta) = KL(\tilde{P}(\mathbf{Z}|\psi), P(\mathbf{Z}|\mathbf{D}; \theta)) \quad (18)$$

となる. $KL(\cdot, \cdot)$ は KL ダイバージェンスであり, 二つの分布の距離を表している. この式より $\tilde{P}(\mathbf{Z}|\psi)$ を最良の近似にするためには $F(\tilde{P}, \theta)$ を最大化するような ψ を用いれば良いことが分かる.

4.2 一般化 EM 法を用いた HMM パラメータの推定

3.2 と同様にシンボル系列 $\mathbf{D} = \{x_1, x_2, \dots, x_K\}$ が観測されたもとで HMM パラメータを GEM によって推定することを考える. いま, 未観測データを $\mathbf{Z} = \{z_1, z_2, \dots, z_{K+1}\}$ とし, その事後分布の近似を $q(\mathbf{z})$ で表現する. つまり,

$$P(\mathbf{Z}|\mathbf{D}) \cong q(\mathbf{Z}). \quad (19)$$

また $q(\mathbf{z})$ に対して, 因数分解が可能である仮定 $q(\mathbf{Z}) = q_1(z_1)q_2(z_2) \dots q_{K+1}(z_{K+1})$ を設けると式 (17) の下限値は,

$$F(q_1, \dots, q_{K+1}; \theta) = \sum_{z_1} \sum_{z_2} \dots \sum_{z_{K+1}} q_1(z_1)q_2(z_2) \dots q_{K+1}(z_{K+1}) \log \frac{P(\mathbf{D}, z_1, z_2, \dots, z_{K+1}; \theta)}{q_1(z_1)q_2(z_2), \dots, q_{K+1}(z_{K+1})} \quad (20)$$

となる. 制約条件 $\sum_{z_k} q_k(z_k) = 1, k = 1, \dots, K+1$ のもとでオイラーラグランジュ方程式を解くと,

$$\begin{aligned} q_k^*(z_k|\mathbf{D}) &\propto \exp\left\{ \sum_{z_1} \dots \sum_{z_{k-1}} \sum_{z_{k+1}} \dots \sum_{z_{K+1}} q_1^*(z_1|\mathbf{D}) \dots q_{k-1}^*(z_{k-1}|\mathbf{D}) q_{k+1}^*(z_{k+1}|\mathbf{D}) \dots q_{K+1}^*(z_{K+1}|\mathbf{D}) \right. \\ &\quad \left. \times \log P(\mathbf{D}, z_1, \dots, z_{K+1}) \right\} \end{aligned} \quad (21)$$

を得る. $q_k^*(z_k)$ は下限値を最大にする $q_k(z_k)$ である. 具体的に $q_k(z_k) = \prod_{i=1}^M \psi_{k,i}^{z_{k,i}}$ とすると

$$\begin{aligned} q_k^*(z_k) \propto & \exp\left\{\sum_{i=1}^M \psi_{1,i} \log \pi_i + \sum_{l=2}^{k-2} \sum_{i=1}^M \sum_{j=1}^M \psi_{l,i} \psi_{l+1,j} \log[B(x_k)A]_{i,j}\right. \\ & + \sum_{i=1}^M \sum_{j=1}^M \psi_{k-1,i} z_{k,j} \log[B(x_k)A]_{i,j} + \sum_{i=1}^M \sum_{j=1}^M z_{k,i} \psi_{k+1,j} \log[B(x_k)A]_{i,j} \\ & \left. + \sum_{l=k+1}^K \sum_{i=1}^M \sum_{j=1}^M \psi_{l,i} \psi_{l+1,j} \log[B(x_k)A]_{i,j} + \sum_{i=1}^M \psi_{K+1,i} \log \xi_i\right\}. \end{aligned} \quad (22)$$

式(22)で $z_{k,i}$ を含む項に注目すると,

$$q_k^*(z_k) \propto \exp\left\{\sum_{i=1}^M \sum_{j=1}^M \psi_{k-1,i} z_{k,j} \log[B(x_k)A]_{i,j} + \sum_{i=1}^M \sum_{j=1}^M z_{k,i} \psi_{k+1,j} \log[B(x_k)A]_{i,j}\right\} \quad (23)$$

となり, $q_k^*(z_k) = \prod_{i=1}^M \psi_{k,i}^{z_{k,i}}$ から

$$\psi_{k,i}^* \propto \prod_{i=1}^M [B(x_k)A]_{i,j}^{\psi_{k-1,i}^*} \times \prod_{i=1}^M [B(x_k)A]_{j,i}^{\psi_{k+1,i}^*} \quad (24)$$

が得られる. 同様に $k=1, K+1$ の時はそれぞれ,

$$\psi_{1,i}^* \propto \pi_i \prod_{j=1}^M [B(x_1)A]_{i,j}^{\psi_{2,j}^*}, \quad \psi_{K+1,i}^* \propto \prod_{j=1}^M [B(x_K)A]_{j,i}^{\psi_{K,j}^*} \xi_i \quad (25)$$

となる. つまり式(24)および式(25)を満たす $\psi_{k,i}^*$ は事後分布 $P(\mathbf{Z}|\mathbf{D})$ に対する最良の近似である. この $\psi_{k,i}^*$ を用いると GEM における E ステップは ψ_k の算出と期待値 $E[N_{i,j}|\mathbf{D}]$, $E[Y_{i,l}|\mathbf{D}]$, $E[R_i|\mathbf{D}]$ の計算となり M ステップはそれらの期待値から式(16)でパラメータを更新する手続きとなる. ψ の算出を組み込んだ GEM によるパラメータ推定は以下ようになる.

Step1 時刻 $k=1, K+1$ および $k=2, 3, \dots, K$ について $\psi_{k,i}$ を更新する.

$$\psi_{1,i} \propto \pi_i \prod_{j=1}^M [B(x_1)A]_{i,j}^{\psi_{2,j}}, \quad (26)$$

$$\psi_{K+1,i} \propto \prod_{j=1}^M [B(x_K)A]_{j,i}^{\psi_{K,j}} \xi_i, \quad (27)$$

$$\psi_{k,i} \propto \prod_{i=1}^M [B(x_k)A]_{i,j}^{\psi_{k-1,i}} \times \prod_{i=1}^M [B(x_k)A]_{j,i}^{\psi_{k+1,i}}. \quad (28)$$

Step2 各期待値を求める.

$$E[N_{i,j}|\mathbf{D}] = \sum_{k=1}^K \psi_{k,i} \psi_{k+1,i}, \quad (29)$$

$$E[Y_{i,l}|\mathbf{D}] = \sum_{k=1}^K I(x_k = v_l) \psi_{k,i}, \quad (30)$$

$$E[R_i|\mathbf{D}] = \psi_{1,i}. \quad (31)$$

Step3 BW と同様に式(16)を用いてパラメータを更新する.

GEM では Step1 から Step3 をパラメータまたは下限値が収束するまで繰り返す.

ここで ψ' はひとつ前の繰り返しで求めた ψ の値であることを示している.

最後に GEM を適用した学習の並列化について考える. 図 2 は GEM を適用した学習の並列化の概念図である. Step1 の ψ^* の更新式に注目すると, 各 $k = 1, \dots, K$ に対する ψ_k の計算では先の反復における ψ に対するデータ依存のみで, $\psi_k, k = 1, \dots, K$ 間のデータ依存がない. よって, スレッド数に応じた並列計算が可能となる. ステップ 2 においても, 和を求めるだけであるから, スレッド数に応じた並列計算が可能である. このように, GEM を適用した学習では効率的に並列化を行うことが可能であり, 長い系列データを扱う際, 大きな効果を生むと考えられる.

5 数値例

5.1 数値実験の環境

ここでは数値実験の環境について記述する. 実験に用いた計算機の CPU は Quad Core AMD Opteron(TM) 2384 (2.7GHz/6MB L3 キャッシュ) を 2 つ用いている. 各 CPU は 4 つのコアを持つので, 原理的に 8 つのプロセスに対する処理を並列に行うことができる. またメモリは 4GB で構成されている.

数値計算用プログラムは, Fortran を使用しており, GEM の GE ステップ, BW の E ステップ, そして GEM と BW 共通の M ステップを実装し, それを C 言語で作成した main 関数によって呼び出している. 実装において, BLAS を用いている. BLAS とは, ベクトルと行列に関する基本線代数数操作を実行するライブラリ API のデファクトスタンダードであり, BLAS によって, GEM と BW の行列演算は最適化されている. 並列化は OpenMP を用いている. OpenMP はメモリ共有型の並列化を行うための API であり, コメント形式の指示文を挿入することでスレッドが生成され, 低コストで並列計算を行えるようになっている. 本稿ではプロセス数を踏まえ, 最大 8 スレッドでの計算を行っている.

数値実験の初期値および教師データは乱数により生成した. また BW と GEM を比較する際の初期値は全て同じ値である. 各数値実験の状態数 n とシンボル数 v は 5, 10, 25, 50, 75, 100 と変化させ, 系列長 $dsize$ は, 10, 100, 1000, 10000 と変化させ数値実験を行った. スレッド数 $thread$ は, 1, 2, 4, 6, 8 と変化させた.

5.2 GEM を適用した HMM の適合性

本稿の以下の数値実験では, HMM の教師データへの適合性を計るために, GEM を適用した HMM の教師データに対する対数尤度の代わりに対数尤度の下限値を用いている. これは, GEM を適用した学習において, 直接対数尤度を求めるためには, BW の一部を流用しなければならず, 収束の判定を行う際, 余計に時間がかかってしまうためである. そこで, ここでは下限値を対数尤度の代用することに対する検証を行う. 下限値は式 (17) から, 完全データの対数尤度と近似分布によって求めることができる. 図 1 は状態数とシ

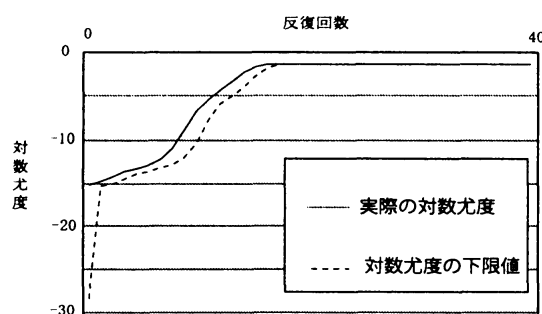


図 3: GEM の対数尤度と下限値の収束.

ンボル数が 10, 教師データの系列長が 10 の時の対数尤度とその下限値の動きを表している. 横軸が反復回数であり, 縦軸が対数尤度である. まず, 下限値が非減少であり, 収束していることがわかる. そして下限値が収束している部分では対数尤度も収束している.

表 1 は状態数, シンボル数, 系列長を変化させながら反復を 1000 回行った時の対数尤度とその下限値である. 表中の値はそれぞれ, 対数尤度とその下限値である. よく似た値を出していることがわかる. 以上のことより, HMM の教師データへの適合性を計る際, 下限値を対数尤度の代わりに近似値として用いることが許容されると言える.

表 1: GEM による学習の対数尤度と下限値.

	dsize=10		dsize=100		dsize=1000		dsize=10000	
	対数尤度	下限値	対数尤度	下限値	対数尤度	下限値	対数尤度	下限値
n=v=5	-6.07	-6.07	-138.95	-139.65	-1591.17	-1591.32	-16086.12	-16087.99
n=v=10	-1.39	-1.39	-168.33	-168.99	-2236.97	-2239.40	-23011.54	-23012.17
n=v=25	0.00	0.00	-142.71	-142.71	-2813.68	-2824.27	-32142.57	-32143.04
n=v=50	0.00	0.00	-95.03	-95.03	-2912.64	-2913.93	-38773.72	-38794.73

5.3 同じ反復回数での実行時間の差

次に GEM による学習と BW が並列化によって計算時間がどのように変化するかを観察する. ここでは計算時間を実行時間と CPU 時間に分けて考える. 実行時間は, 実際に計算が開始されてから終了するまでの

表 2: HMM 学習の実行時間.

n=v=50		dsize=10									
		thread=1		thread=2		thread=4		thread=6		thread=8	
		実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間
BW		0.71	0.71	0.56	1.13	0.52	2.09	0.56	3.35	0.59	4.69
GEM		1.18	1.17	0.96	1.92	0.89	3.55	0.92	5.51	0.93	7.39

n=v=50		dsize=10000									
		thread=1		thread=2		thread=4		thread=6		thread=8	
		実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間	実行時間	CPU 時間
BW		1405.75	1405.47	731.04	1438.45	471.53	1521.19	835.40	1607.81	341.69	1687.30
GEM		603.30	603.18	301.11	602.11	150.33	601.21	100.98	605.74	76.04	608.21

時間である. CPU 時間は全てのコアでのクロック数の合計を, 一秒あたりのクロック数で割って計算している. どちらも単位は秒である.

表 2 は状態数及びシンボル数を 50 に固定した時の, 実行時間と CPU 時間を表わしている. GEM と BW を短い系列で比較すると, どちらもスレッドの増加に従って実行時間が短縮されるものの, 系列長が 10 の場合では, スレッド数が 6 から 8 になった時に, 実行時間の値が増加している. また, 系列長が 100 の場合でもスレッドの増加による実行時間の減少幅も小さい. この時の CPU 時間を比較すると, BW, GEM とともに, スレッドの増加とともに CPU 時間が増加している. CPU 時間は CPU の仕事量と捉えることもできる. スレッドが増加するにしたがって, CPU 時間が増加しているのはスレッド処理のためのオーバーヘッドが増加しているためと考えられる. 短い系列の学習において CPU 時間の増加が大きく, 実行時間の短縮幅が小さいのは, オーバーヘッドが占める割合が大きいためである.

次に長い系列の場合を比較する. 系列長が 10000 の場合, スレッド数が 1 から 8 に変化すると, 実行時間は, BW で 4 分の 1, GEM では 8 分の 1 程度になっている. これは, 実行時間に占めるスレッド処理のオーバーヘッドの割合が小さくなり並列化の効果が顕著に現れることを示している. また, BW はスレッドの増加に対して, 計算時間の短縮幅は頭打ちになっているように見えるが, GEM ではその傾向は見られない. アムダールの法則により, スレッド数をどんどん増やしていくと, GEM でも実行時間の短縮幅は頭打ちになると考えられるが, 8 並列まででは, ほぼスレッド数と実行時間が反比例している. これは, BW では E ステップで前向き確率と後ろ向き確率の計算において, データ依存を含んでおり, 高々 2 並列での計算しか行うことができないのに対し, GEM では, 各時刻ごとに並列化が可能であるためだと考えられる. これらのことから, 長い系列を計算する場合, GEM において, 並列化効果が特に顕著に現れており, GEM は BW より効率的に並列計算が行えると言える.

5.4 収束速度の比較

次に, 収束速度を比較する. またその時の BW と GEM を適用した HMM の教師データに対する, 対数尤度, あるいは対数尤度の下限値を見る. これにより GEM が HMM の学習として利用できるかどうかを考察する. 実際の計算において, 回数で打ち切るよりも, 収束の判定によって打ち切ることが多い. 収束の判

表 3: BW と GEM による学習の収束速度の比較.

n=v=50		dsize=100					
	反復回数	対数尤度	thread=1	thread=2	thread=4	thread=6	thread=8
BW	34	-97.22	0.24	0.14	0.09	0.08	0.07
GEM	40	-95.03	0.26	0.14	0.09	0.07	0.07

n=v=50		dsize=10000					
	反復回数	対数尤度	thread=1	thread=2	thread=4	thread=6	thread=8
BW	3202	-36945.78	3830.44	1972.80	1277.28	1047.38	931.44
GEM	2423	-37643.74	1426.17	713.52	357.73	238.50	179.40

定は相対誤差と絶対誤差がそれぞれの基準値より下回った時に、収束したと判定した。今回は相対誤差の基準値は 10^{-7} 、絶対誤差の基準値は 10^{-5} としている。表 3 は BW と GEM による学習の収束速度を比較したものである。状態数とシンボル数が 50 の場合の系列長が 100 の時と、10000 の時を示している。反復回数は収束するまでに要した反復の回数であり、それ以外の数値は実行時間を示している。収束するまでの反復回数は BW が多ければ GEM も多く BW が少なければ GEM も少ない傾向がある。計算時間は上で示したように GEM の方が顕著に効果が現れている。GEM と BW の対数尤度には開きがある。これは BW、GEM いずれも収束しきっていないためだと考えられる。収束の判定条件を厳しくすることで、より近い値が出せると考えられる。

6 まとめと今後の課題

HMM に一般化 EM 法を適用することで並列化に適した学習アルゴリズムを構成できることが分かった。並列化を行うことにより GEM は BW より効率的な並列化を実現できる。しかし系列が短い場合は、その効果が現れにくく、その精度は BW に劣るという結論を導ける。また対数尤度の比較からある程度有効な学習が行われていることも分かった。実際の適用において BW を適用した HMM と GEM を適用した HMM でどの程度性能に差があるかを、出力確率に確率密度関数を持つ連続型 HMM での検証やトラフィック解析への適用を通して調べて行く予定である。

参考文献

- [1] 中川 聖一, 「確率モデルによる音声認識」, 電子情報通信学会, 1988.
- [2] 前村 一哉, 小野 廣隆, 定兼 邦彦, 山下 雅史, 隠れマルコフモデルを用いた DNA 配列設計, 電子情報通信学会技術報告, vol. 106, no. 63, pp. 39-46, 2006.
- [3] 菊池 洋光, 甲藤 二郎, 改良型 Embedded HMM を用いた顔画像認識の検討, 画像電子学会研究会講演予稿, Vol. 02-07, pp. 85-90, 2003.
- [4] L.E. Baum, T. Petrie, G. Soules and N. Weiss, A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains, *Ann. Math. Statist.*, vol. 41, no. 1, pp. 164-171, 1970.
- [5] 松野 裕之, 友利 記昌, 宮崎 崇, 西村 英樹, 神戸 尚志, 音声認識システムのハードウェア化の一手法—HMM 出力確率計算のハードウェア化, 電子情報通信学会技術研究報告, vol. 104, no. 737, pp. 79-84, 2005.
- [6] 汪 金芳, 田栗 正章, 手塚 集, 樺島 洋介, 上田 修功, 「計算統計 I 確率計算の新しい手法」, 岩波書店, 2003.
- [7] Z. Ghahramani and M.I. Jordan, Factorial hidden markov models, *Machine Learning*, vol. 29, pp. 245-275, 1997.