



東京大学 工学部 計数工学科/物理工学科

応用音響学 : Baum-Welch アルゴリズム

嵯峨山 茂樹 <sagayama@hil.t.u-tokyo.ac.jp>

東京大学 工学部 計数工学科 <http://hil.t.u-tokyo.ac.jp/>



EM アルゴリズムと Baum-Welch アルゴリズム

■ 内容

- HMM の 3 つの問題
- EM アルゴリズム
- EM アルゴリズム収束性の証明
- Baum-Weich アルゴリズム

■ 参考文献

- 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.
- Lawrence Rabiner 他「音声認識の基礎 (下)」NTTアドバンステクノロジー株式会社, 1995.
- 上坂吉則他「パターン認識と学習のアルゴリズム」文一総合出版, 1990.
- 北研二他「音声言語処理 コーパスに基づくアプローチ」森北出版株式会社, 1996.



HMMの3つの基本問題

モデル $\lambda = \left(\begin{array}{c} \text{状態遷移確率行列} \\ A \end{array}, \begin{array}{c} \text{観測シンボル確率分布} \\ B \end{array}, \begin{array}{c} \text{初期状態分布} \\ \pi \end{array} \right)$ 、
観測系列 $Y = (y_1 y_2 \dots y_T)$

問題1 モデル λ に対する観測系列 Y の確率 $P(Y|\lambda)$ の計算

- モデル λ が観測系列 Y に対してどの程度適応しているか
... ex. Forward アルゴリズム

問題2 最適な状態系列 $q = (q_1 q_2 \dots q_T)$ の発見

- 観測系列 Y がどの状態系列 q から生成されたと考えられるか
... ex. Viterbi アルゴリズム

問題3 $P(Y|\lambda)$ を最大とするようなモデルパラメータ λ の調整

- 観測系列 Y を生成するためのパラメータ λ の最適化
... ex. EM アルゴリズム、Baum-Welch アルゴリズム

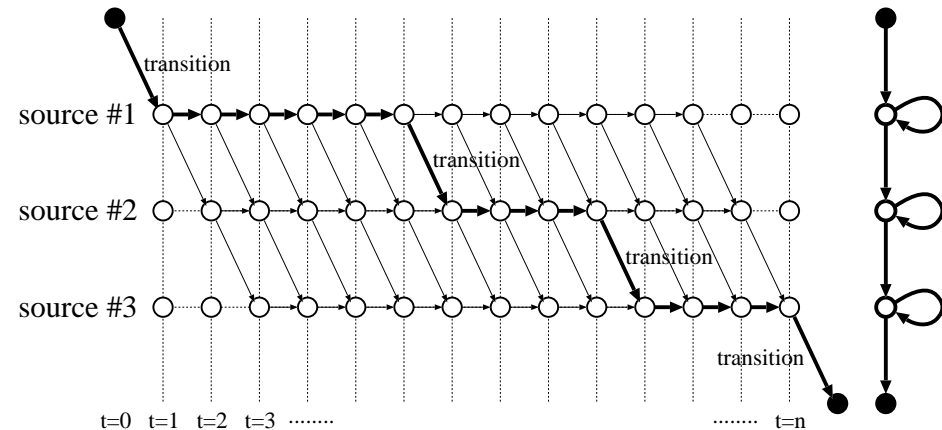


Viterbi 学習と Baum-Welch 学習

■ パラメータ学習アルゴリズム

■ Viterbi 学習アルゴリズム:
 確率最大の経路に沿って再学習の繰り返し

■ Baum-Welch 学習アルゴリズム:
 全経路を考えてトレリス各点で確率重みつき再学習の繰り返し



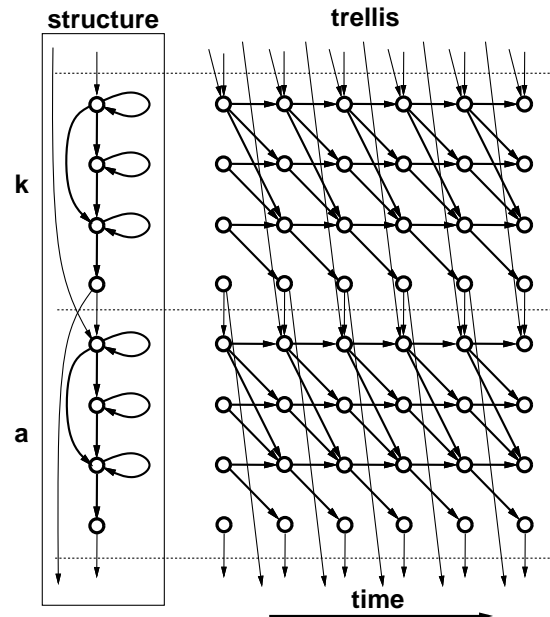


HMMの学習

学習アルゴリズム

- Baum-Welch アルゴリズム (forward-backward algorithm)
- Viterbi 学習アルゴリズム

HMMトレリス計算(縦方向は状態、横方向は時間)



Trellis Computation

図1. HMMトレリス計算(縦方向は状態、横方向は時間)



EM(Expectation-Maximization) アルゴリズム

観測できるデータ x : 不完全データ
観測できないデータ y } 完全データ

完全なデータが得られる場合は最尤推定によってパラメータを推定することができるが、不完全なデータしか得られない場合は最尤推定を直接行なうことはできない。

EM アルゴリズムは、不完全データからの対数尤度を最大化するために、完全データからの対数尤度の**期待値を最大化**する。



EM アルゴリズム

y : 不完全な観測データ、 x : y に付随して得られる付加データ、 ϕ : 初期パラメータ、 $f(y|\theta)$, $f(x|\theta)$: 確率密度関数 とする。

尤度関数 $P(\hat{\theta}; y)$ に x を導入

$$P(\hat{\theta}; y) = \frac{P(\hat{\theta}; x, y)}{P(\hat{\theta}; x, y)P(\hat{\theta}; y)} = \frac{P(\hat{\theta}; x, y)}{P(\hat{\theta}; x|y)} \quad (1)$$

対数尤度

$$\log P(\hat{\theta}; y) = \log P(\hat{\theta}; x, y) - \log P(\hat{\theta}; x|y) \quad (2)$$



EM アルゴリズム

確率変数 X のすべてに渡って、あらかじめ設定されている推定値 θ を使って期待値をとる。

$$E[\log P(\hat{\theta}; y)|x] = E[\sum_x P(\theta, x|y) \log P(\hat{\theta}; y)] \quad (3)$$

$$= \log P(\hat{\theta}; y) \quad (4)$$

$$E[\log P(\hat{\theta}; y)] = E[\log P(\hat{\theta}; x, y) - \log P(\hat{\theta}; x|y)] \quad (5)$$

$$= E[\log P(\hat{\theta}; x, y)] - E[\log P(\hat{\theta}, x|y)] \quad (6)$$

これから

$$\log P(\hat{\theta}; y) = E[\log P(\hat{\theta}; x, y)] - E[\log P(\hat{\theta}; x|y)] \quad (7)$$

が得られる。



EM アルゴリズム

ここで、

$$\log P(\hat{\theta}; y) = E[\log P(\hat{\theta}; x, y)] - E[\log P(\hat{\theta}; x|y)] \quad (8)$$

$$= \sum_x P(\theta; x, y) \log P(\hat{\theta}; x, y) - \sum_x P(\theta; x|y) \log P(\hat{\theta}; x|y) \quad (9)$$

$$= Q(\theta, \hat{\theta}) - H(\theta, \hat{\theta}) \quad (10)$$

Jensen の不等式より、

$H(\theta, \hat{\theta}) \leq H(\theta, \theta)$ なので、 $Q(\theta, \hat{\theta}) \geq Q(\theta, \theta)$ となるように $\hat{\theta}$ を設定すれば $P(y|\hat{\theta}) \geq P(y|\theta)$ となる。

つまり、 **Q を最大化するよう θ を変えると P も最大化される。**

EM アルゴリズムは $P(y|\theta)$ を直接最大化する代わりに Q を最大化していくアルゴリズムである。



EM アルゴリズム

Jensen の不等式

$$H(\Phi, \Phi) \geq H(\Phi, \bar{\Phi})$$

$$H(\Phi, \Phi) = \int \log f(y|x, \Phi) f(y|x, \Phi) dy$$

$$H(\Phi, \bar{\Phi}) = \int \log f(y|x, \bar{\Phi}) f(y|x, \Phi) dy$$

$$\begin{aligned} H(\Phi, \bar{\Phi}) - H(\Phi, \Phi) &= \int (\log f(y|x, \bar{\Phi}) - \log f(y|x, \Phi)) f(y|x, \Phi) dy \\ &= \int \log \frac{f(y|x, \bar{\Phi})}{f(y|x, \Phi)} f(y|x, \Phi) dy \\ &\leq \int \left(\frac{f(y|x, \bar{\Phi})}{f(y|x, \Phi)} - 1 \right) f(y|x, \Phi) dy \quad \text{since } \log A \leq A - 1 \\ &= \int f(y|x, \bar{\Phi}) dy - \int f(y|x, \Phi) dy = 0 \end{aligned}$$

等号は $f(y|x, \bar{\Phi}) = f(y|x, \Phi)$ の時成立



EM アルゴリズム

これを観測可能な系列 $X = \{x_1, \dots, x_n\}$ と観測不可能な系列 $Y = \{y_1, \dots, y_n\}$ の場合に拡張すると以下のようなになる。

$$L(X, \theta) = \sum_{k=1}^n \log f(x_k | \theta) \quad (11)$$

$$E[L(X, \hat{\theta}) | X, \theta] = E\left[\sum_{k=1}^n \log f(x_k | \theta) | X, \theta\right] \quad (12)$$

$$= L(X, \theta) \quad (13)$$

$$Q(\theta, \hat{\theta}) = \sum_{k=1}^n Q_k(\theta, \hat{\theta}) \quad (14)$$

$$H(\theta, \hat{\theta}) = \sum_{k=1}^n H_k(\theta, \hat{\theta}) \quad (15)$$



EM アルゴリズム

実行手順

1. パラメータ Φ の初期値を設定
2. $Q(\Phi, \bar{\Phi})$ を求める。
3. $Q(\Phi, \bar{\Phi})$ を最大にするような $\bar{\Phi}$ を選ぶ。
4. $\bar{\Phi}$ を Φ に設定し、収束条件が満たされなければ 2. へ、満たされれば終了

ステップ2. は**期待値操作 (Expectation step)**、ステップ3. は**最大値操作 (Maximization step)**と呼ばれている。



Baum-Welch アルゴリズム

- 観測系列の確率を最大化するモデルパラメータを解析的に直接求める方法は知られていない
- しかし尤度が局所的に最大になるモデルパラメータを求めることはできる
- EM アルゴリズムをHMMのパラメータ推定に適用... **Baum-Welch アルゴリズム (Forward-Backward アルゴリズム)**

パラメータの推定

モデルと観測系列が与えられたとき、時刻 t のとき状態 i に存在し、時刻 $t + 1$ のとき状態 j に存在する確率 $\xi_t(i, j)$ を定義する。

$$\xi_t(i, j) = P(q_t = i, q_{t+1} = j | O, \lambda)$$

上式の条件を満足するパスを以下に示す。



Baum-Welch アルゴリズム

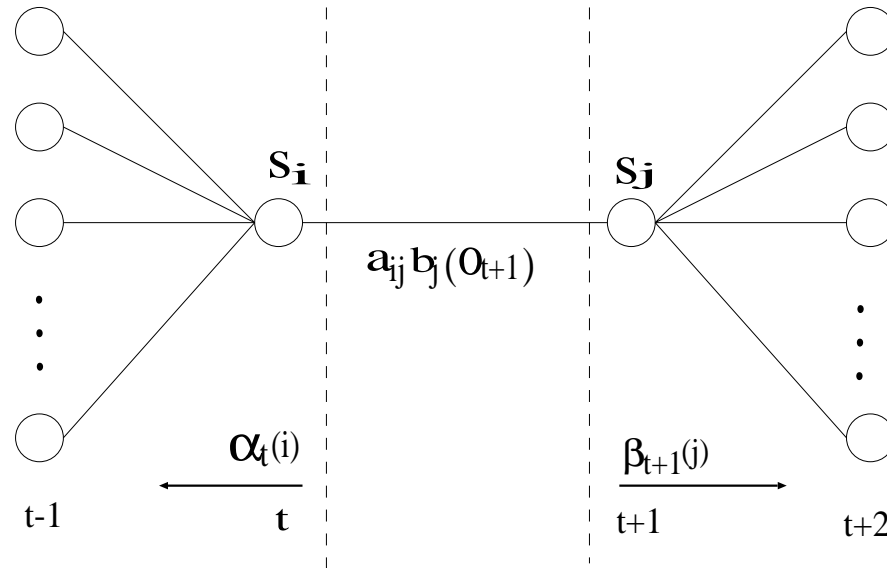


図2. 計算手順

前向き変数を α 、後向き変数を β とすると $\xi_t(i, j)$ は次のように書ける。

$$\begin{aligned} \xi_t(i, j) &= \frac{P(q_t = i, q_{t+1} = j, Y | \lambda)}{P(Y | \lambda)} \\ &= \frac{\alpha_t(i) a_{ij} b_j(y_{t+1}) \beta_{t+1}(j)}{P(Y | \lambda)} \end{aligned}$$



Baum-Welch アルゴリズム

$$= \frac{\alpha_t(i) a_{ij} b_j(y_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(y_{t+1}) \beta_{t+1}(j)}$$

$\gamma_t(i)$ をモデルと観測系列全体が与えられたときに時刻 t で状態 i に存在する確率と定義すると、 $\gamma_t(i)$ は $\xi_t(i, j)$ を j について総和したものと考えられるので、

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

$\gamma_t(i)$ を t について総和をとると状態 i を訪れた回数の期待値とみなせ、それから $t = T$ を除いた値は i から遷移する回数の期待値とみなせる。

$$\sum_{t=1}^{T-1} \gamma_t(i) : Y \text{ において状態 } i \text{ から遷移する回数の期待値}$$

同様に

$$\sum_{t=1}^{T-1} \xi_t(i, j) : Y \text{ において状態 } i \text{ から状態 } j \text{ へ遷移する回数の期待値}$$

とみなせる。



Baum-Welch アルゴリズム

これらにより、HMMのパラメータ π (初期状態分布), A (状態遷移確率行列), B (観測シンボル確率分布) の再推定式は以下ようになる。

$$\overline{\pi}_i = \text{時刻 } (t = 1) \text{ に状態 } i \text{ に存在すると期待される頻度 (回数)} \quad (16)$$

$$= \gamma_1(i) \quad (17)$$

$$\overline{a}_{ij} = \frac{\text{状態 } i \text{ から状態 } j \text{ へ遷移する回数の期待値}}{\text{状態 } i \text{ から遷移する回数の期待値}} \quad (18)$$

$$= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \quad (19)$$

$$\overline{b}_j(k) = \frac{\text{状態 } j \text{ にとどまりシンボル } v_k \text{ を観測する回数の期待値}}{\text{状態 } j \text{ にとどまる回数の期待値}} \quad (20)$$

$$= \frac{\sum_{t=1}^T \sum_{s.t. o_t=v_k} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad (21)$$



Baum-Welch アルゴリズム

現在のモデルを $\lambda = (A, B, \pi)$ 、(17),(18),(20) 式の左辺によって決定される再推定モデルを $\bar{\lambda} = (\bar{A}, \bar{B}, \bar{\pi})$ とすると、

1. 初期モデルが尤度関数の臨界点 ($\bar{\lambda} = \lambda$ の点) を定義する、あるいは
2. $\bar{\lambda}$ が λ よりも $P(Y|\bar{\lambda}) > P(Y|\lambda)$ の意味でより尤もらしい

つまり、観測系列が生成された可能性がより高い新しいモデル $\bar{\lambda}$ を手に入れることができる。

このように $\bar{\lambda}$ を λ に入れ換えて繰り返し使いながら再推定計算を繰り返せば Y がそのモデルから観測されたという確率をある点まで高めることができるが、これはあくまでも極大点である。



Baum-Welch アルゴリズム

(17),(19),(21)式の再推定式は次の補助関数を最大化することで直接求めることができる。

$$Q(\lambda', \lambda) = \sum_q P(Y, q|\lambda') \log P(Y, q|\lambda)$$

なぜなら

$$Q(\lambda', \lambda) \geq Q(\lambda', \lambda') \Rightarrow P(Y|\lambda) \geq P(Y|\lambda')$$

の関係があるので $P(Y|\lambda)$ を増加させるという意味で λ に関する関数 $Q(\lambda', \lambda)$ を λ' を改善しながら最大化できるからである。これを繰り返せば尤度関数は最終的に臨界点に達する。



Baum-Welch アルゴリズム

再推定式の導出

P と $\log P$ は HMM パラメータにより次のように表現できる。

$$P(Y, q|\lambda) = \pi_{q_0} \prod_{t=1}^T a_{q_{t-1}q_t} b_{q_t}(y_t)$$

$$\log P(Y, q|\lambda) = \log \pi_{q_0} + \sum_{t=1}^T \log a_{q_{t-1}q_t} + \sum_{t=1}^T \log b_{q_t}(y_t)$$

よって

$$Q(\lambda', \lambda) = Q_{\pi}(\lambda', \pi) + \sum_{i=1}^N Q_{a_i}(\lambda', a_i) + \sum_{i=1}^N Q_{b_i}(\lambda', b_i)$$

そして

$$Q_{\pi}(\lambda', \pi) = \sum_{i=1}^N P(Y, q_0 = i|\lambda') \log \pi_i$$

$$Q_{a_i}(\lambda', a_i) = \sum_{j=1}^N \sum_{t=1}^T P(Y, q_{t-1} = i, q_t = j|\lambda') \log a_{ij}$$

$$Q_{b_i}(\lambda', b_i) = \sum_{t=1}^T P(Y, q_t = i|\lambda') \log b_i(y_t)$$



Baum-Welch アルゴリズム

$Q(\lambda', \lambda)$ は3つの独立な項に分かれているので次の制約の下でそれぞれの項を最大化することで $Q(\lambda', \lambda)$ を λ に関して最大化できる。

$$\sum_{j=1}^N \pi_j = 1$$

$$\sum_{j=1}^N a_{ij} = 1, \forall i$$

$$\sum_{k=1}^K b_i(k) = 1, \forall i$$

これらは全て

$$\sum_{j=1}^N w_j \log y_j$$

の形をしている。これは $\sum_{j=1}^N y_j = 1, (y_j \geq 0)$ の制約が存在する $\{y_j\}_{j=1}^N$ の関数として次のような極大値に至る。

$$y_j = \frac{w_j}{\sum_{i=1}^N w_i}, \quad j = 1, 2, \dots, N$$



Baum-Welch アルゴリズム

以上より、最大化処理は再推定モデル $\bar{\lambda} = [\bar{\pi}, \bar{A}, \bar{B}]$ を次の点に導く。

$$\bar{\pi}_i = \frac{P(Y, q_0 = i | \lambda)}{P(Y | \lambda)}$$

$$\bar{a}_{ij} = \frac{\sum_{t=1}^T P(Y, q_{t-1} = i, q_t = j | \lambda)}{\sum_{t=1}^T P(Y, q_{t-1} = i | \lambda)}$$

$$\bar{b}_i(k) = \frac{\sum_{t=1}^T P(Y, q_t = i | \lambda) \delta(y_t, v_k)}{\sum_{t=1}^T P(Y, q_t = i | \lambda)}$$

ここで、

$$\delta(y_t, v_k) = \begin{cases} 1 & (y_t = v_k) \\ 0 & (\text{otherwise}) \end{cases}$$



Baum-Welch アルゴリズム

前向き変数 α と後向き変数 β を使うと

$$P(Y, q_t = i | \lambda) = \alpha_t(i) \beta_t(i)$$

$$P(Y | \lambda) = \sum_{i=1}^N \alpha_t(i) \beta_t(i) = \sum_{i=1}^N \alpha_T(i)$$

$$P(Y, q_{t-1} = i, q_t = j | \lambda) = \alpha_{t-1}(i) a_{ij} b_j(y_t) \beta_t(j)$$

これらより、前ページの式はそれぞれ**(2)**,**(4)**,**(6)**式であらわされる。

**証明**

$\log Z \leq Z - 1$ の関係より、

$$\begin{aligned} Q(\lambda', \lambda) - Q(\lambda', \lambda') &= \sum_q P(Y, q|\lambda') \log P(Y, q|\lambda) - \sum_q P(Y, q|\lambda') \log P(Y, q|\lambda') \\ &= \sum_q P(Y, q|\lambda') \log \frac{P(Y, q|\lambda)}{P(Y, q|\lambda')} \\ &\leq \sum_q P(Y, q|\lambda') \left\{ \frac{P(Y, q|\lambda)}{P(Y, q|\lambda')} - 1 \right\} \\ &= \sum_q \{P(Y, q|\lambda) - P(Y, q|\lambda')\} \\ &= P(Y|\lambda) - P(Y|\lambda') \end{aligned}$$

以上より

$$Q(\lambda', \lambda) \geq Q(\lambda', \lambda') \Rightarrow P(Y|\lambda) \geq P(Y|\lambda')$$