

効率的なクロスバリデーションに基づく混合ガウス分布の 最適化法とその拡張

篠崎 隆宏[†] 古井 貞熙[†] 河原 達也[‡]

[†] 東京工業大学 大学院情報理工学研究科 計算工学専攻

[‡] 京都大学 学術情報メディアセンター

あらまし 有限の学習データからの混合ガウス分布の推定において、混合要素の数およびそれらの配置を最適化することは、学習データに含まれない新しいデータに対し高い性能を得る上で非常に重要である。これまでに混合ガウス分布を最適化する手法としてクロスバリデーション (CV) を効率的に適用する手法の提案を行いその有効性を示したが、本研究ではモデル性能の更なる向上を目的として、CV 手法を拡張した Aggregated CV (AgCV) 法および AgCV 法を混合ガウス分布の最適化に応用する手法の提案を行なう。提案アルゴリズムは従来の CV を用いた混合ガウス分布の最適化法と同様に、十分統計量を用いることで効率的に動作する。日本語話し言葉コーパスを用いた大語彙連続音声認識実験において、混合ガウス分布 HMM の最適化に本手法を用いることで、モデルサイズを自動決定しつつ従来法と比較して認識性能が向上することを示す。

Gaussian Mixture Optimization Based on Efficient Cross-validation and Its Extension

Takahiro SHINOZAKI[†] Sadaoki FURUI[†] Tatsuya KAWAHARA[‡]

[†] Tokyo Institute of Technology

[‡] Kyoto University

Abstract We have previously proposed a cross-validation (CV) based Gaussian mixture optimization method that efficiently optimizes the model structure based on CV likelihood. In this study, we first propose aggregated cross-validation (AgCV) that extends CV, and then apply it to Gaussian mixture optimization to further improve the model performance. The AgCV based Gaussian mixture optimization algorithm works efficiently by utilizing sufficient statistics similarly to the CV based optimization method. The proposed algorithm is applied to Gaussian mixture HMM and evaluated by speech recognition experiments on oral presentations. It is shown that lower word error rates than conventional methods are obtained by the AgCV optimization method incorporating automatic model size determination.

1 はじめに

混合ガウス分布は混合ガウス分布モデル (GMM) や混合ガウス分布 HMM などとして話者認識や音声認識など幅広く用いられている。混合ガウス分布の推定において一般に問題となるのが、如何にしてモデル性能が最大となるように、与えられた有限の学習データに対してパラメタ数 (混合要素数) とパラメタの推定精度のバランスをとるかという問題である。また、混合ガウス分布は混合重みとして隠れ変数を持ち学習において多くの局所最適解を持つことから、混合要素の数とともにそれらをどのように配置するかということも、高い性能を得る上で重要となる。

大きな混合数を持つ混合ガウス分布が与えられた場合、混合要素の配置および混合数を最適化するひとつの戦略として、適当な目的関数を基準として、混合要素対の選択・併合を停止条件が満たされるまで繰り返していくことが考えられる。この場合混合要素対を一つ併合する毎に全ての混合要素の組に対して目的関数を評価する必要があり、目的関数が効率的に計算されることが実用上で非常に重要となる。

最も一般的な目的関数としては、モデルの学習セットに対する尤度が考えられる。尤度を混合ガウス分布最適化の目的関数とすることは、尤度を用いたパラメタ学習と一貫性があるとともに、十分統計量を用いることで高速に評価可能であるという利点がある [1]。しかしながら、学習セットに対する尤度を目的関数とした場合、評価値が要素対併合に対し常に減少し、モデルパラメタ数とモデル推定精度のバランスをとる停止基準が与えられない問題がある。これは、モデルのパラメタ推定と尤度の評価が同一の学習セットに対して行われるため、尤度が正のバイアスを含むためである。対処法として適当な閾値を用いることも考えられるが、閾値を経験的に調節することが必要となる。また、情報量基準を用いることで停止基準を与えることも出来るが、実際には理論値からの誤差を補正するための経験的な補正項がしばしば必要となる [2]。

クロスバリデーション (CV) はモデルの推定と評価に用いるデータを効果的に分離することで、モデル評価におけるこのようなバイアスを大きく減らすことが出来るデータ駆動に基づく手法である。これにより評価値の最大値を与える点として、モデルパラメタ数を自動的に決定することも可能となる。

従来 CV を用いたモデル選択は計算量の問題から音声認識に用いられるような大規模なモデルに対する応用においては比較的少数個のモデルの比較や半連続 HMM への適用などに限られていた。このような背景のもと、これまでの研究で逐次状態分割法 [3]

や選択学習法 [4] で用いられている十分統計量を用いる手法と似た方法でガウス分布に対する CV 尤度が効率的に評価可能であり、HMM の状態クラスタリング [5] や混合ガウス分布の最適化 [6] への応用において CV 尤度による最適化手法が有効であることを示した。

しかしながら、CV をこのような構造最適化に用いる場合、従来の CV の利用と比べて比較対象となるモデル数が非常に大きくなることが注意点として挙げられる。たとえば混合ガウス分布の最適化の際には、 M 混合の混合ガウス分布の要素数を 1 つ減らし $M-1$ とするのにも、 $\frac{M(M-1)}{2}$ 個の数のモデルが比較され、さらにこれが停止条件が満たされるまで繰り返される。他方、CV を用いることで評価値バイアスは大部分取り除かれるが、CV による評価値は依然として学習データの分布や CV 区画の切り方などに依存した統計的な揺らぎを持つ。このため、このように非常に多くのモデルを対象としたモデル選択が再帰的に繰り返される場合、評価値の揺らぎに起因して実際の性能を反映しない形で高い評価値を与えるモデルが選択されてしまう可能性が無視できなくなる恐れがある。また、この効果は選択対象モデルの数とともに増大し、最終的にモデル性能の低下となって表れると考えられる。

本研究では、このような評価値の揺らぎを減らすことを目的として Bagging [7] に似たアイデアを CV の枠組内部に導入する Aggregated CV (AgCV) 法の提案を行い、さらに CV の代わりに AgCV を用いた混合ガウス分布の最適化手法を提案する。AgCV 法は Bagging に似たアイデアを用いるという点ではこれまでに提案した AgEM アルゴリズム [8, 9] と似ているが、AgEM が EM [10] を拡張したパラメタ推定アルゴリズムであるのに対し、本論文で提案する AgCV 法は CV を拡張したモデル構造選択アルゴリズムであるという点が、大きく異なる。なお以下では、従来の学習セットに対する尤度を提案法による尤度と区別するため、自己尤度と呼ぶ。

本論文の構成は以下のとおりである。まず CV の拡張である AgCV 法について、第 2 章で提案する。ついで AgCV を用いた混合分布の最適化手法を第 3 章で提案する。第 4 章で実験条件について説明し、第 5 章で実験結果を示す。最後に第 6 章でまとめと今後の課題を示す。

2 CV および Aggregated CV (AgCV) 法

本章では混合ガウス分布の最適化とは独立に、一般的な CV の拡張法としての Aggregated CV (AgCV)

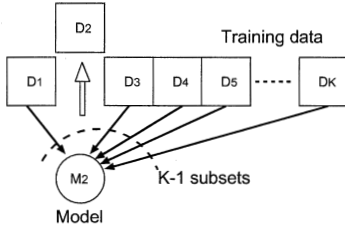


図 1: *K-fold cross-validation (K-fold CV)*.

法の提案を行う。はじめに従来の *K-fold CV* について簡単に説明し、ついで提案法である *AgCV* 法を示す。

2.1 *K-fold* クロスバリデーション (*K-fold CV*)

一般に、*K-fold CV* では学習データを K 個の区画に区分化する。そして、図 1 に示すように、 K 個中 1 つの区画を取り除いた残りの $K - 1$ 区画からモデル推定を行い、得られたモデルを取り除いておいた区画に適用し、評価値の計算を行う。この操作を取り除く区画を K 通りかえて繰り返し、それらの評価値を足し合わせたのが *CV* 評価値である。*CV* 法では学習データをモデルの推定用と評価用に分割するが、各 *CV* モデルは学習データ全体の $(K - 1) / K$ を使用して推定されるため、データの断片化問題は K を大きくとることでほぼ無視できる程度に抑えることができる。*CV* 法を用いることで、モデルの推定と評価におけるデータの重複をなくすことができ、モデル推定と評価に同じデータセットを用いる場合と比較して、大幅に評価値に含まれる正のバイアスを減らすことができる。このため、*CV* 評価値を基準とすることで、学習データに含まれない新しいデータに対するモデル性能を精度良く近似することができ、学習データのみを用いたモデル選択が可能となる。

2.2 *Aggregated CV (AgCV)*

Aggregated CV (AgCV) 法では、*CV* 法による評価値の揺らぎを減らしモデル選択能力を向上させることを目的として、*Bagging* [7] に似たアイデアを *K-fold CV* の枠組の内部に導入する。*Bagging* は複数の分類器の出力を統合することで識別性能を向上させるアンサンブル学習の一手法であり、与えられた学習セットから無作為復元抽出により重複を許した複数のサブ学習セットを定義し、それら各サブセット毎に識別器の学習を行う。そして、識別時にはそ

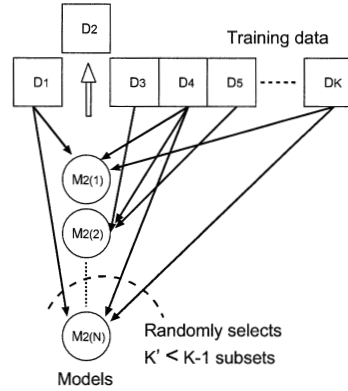


図 2: *Aggregated cross-validation (AgCV)*.

れらの識別器を並列に動作させ、投票や平均操作により最終的な出力を与える。

本研究で提案する *AgCV* 法は *K-fold CV* 法の拡張であり、図 2 に示すように、評価対象となる区画をその区画を含めず推定した N 個のモデルにより繰り返し評価し、それらの平均をとる。評価に用いる N 個のモデルは評価対象の区画を除いた残りのデータから、*Bagging* と同様に相互に重複を許した複数のサブセットを定義し学習することにより得る。ただし本来の *Bagging* とは異なり、 N 個の学習サブセットは $K - 1$ 個の *CV* 区画から $K' (\leq K - 1)$ 個の区画を非復元抽出する“粗い”抽出操作を N 回繰り返すことにより行う。この粗粒度の抽出方法は、*AgCV* を十分統計量を用いたモデル選択に応用する際に、有用となる。もし $N = 1, K' = K - 1$ とした場合には、*AgCV* は従来の *K-fold CV* に一致する。

各区画の評価に用いる N 個のモデル相互の類似度は、モデル間で共有されるデータの割合 K' / K により制御することができる。 K' が大きいと N 個のモデルはどれも似たものとなりアンサンブル評価の効果が期待できず、また小さすぎると個々のモデルの学習に実際に使われるデータ量が少なくなってしまう。本研究では事前実験により、 $K' = \frac{1}{2}K$ としている。

3 *AgCV* 法による混合ガウス分布の最適化

AgCV 法による混合ガウス分布の最適化は、*CV* 法による最適化法の拡張として得られる。本節では、自己尤度法、*CV* 尤度法および *AgCV* 尤度法による最適化アルゴリズムを並行して説明する。どのアル

ゴリズムも大枠は同じであり、異なるのは尤度評価法の計算法のみである。

3.1 モデル最適化手順

提案する混合ガウス分布の最適化法では、入力として大きな混合数を持つモデルを受け取り、過剰な要素を順次削減していく。余剰要素の削減は目的関数に基づきモデルの汎化性能向上の上でもっとも効果的となるように、混合要素の対を選択し併合することにより行う。即ち、 M 混合の混合ガウス分布のどれか 1 対の要素分布を併合することにより得られる $\frac{M(M-1)}{2}$ 個の可能な $M-1$ 混合モデル全ての評価値を計算し、最大値を与えるモデルを選択する。この操作を順次繰り返す、停止基準が満たされるまで混合要素の削減を進める。

目的関数として学習セットに対する自己尤度を用いるのが従来の自己尤度法であり、CV 尤度を用いるのが前回提案した CV 尤度による最適化法、AgCV 尤度を用いるのが今回新たに提案する手法である。自己尤度を用いる場合は、尤度が要素分布の併合に対して単調となるため、停止基準としては尤度の変化量に対して適当な閾値を設定するか、情報量基準を用いる必要がある。CV 尤度、および AgCV 尤度を用いる場合は尤度は単調減少とはならず、これらの尤度の最大値を与える点として、最適な混合数を容易に判断することができる。

いずれの評価尺度を用いる場合においても、非常に多数のモデルを対象とした評価値計算が必要となることから、効率的なアルゴリズムが不可欠となる。以下では十分統計量を用いた、自己尤度、CV 尤度、および AgCV 尤度の効率的な評価アルゴリズムについて説明する。混合ガウス分布 HMM の最適化は、各 HMM 状態ごとに混合ガウス分布の最適化を行うことで得られる。

3.2 十分統計量を用いた評価値の計算

従来の自己尤度法、CV 尤度法、および AgCV 尤度法における効率的な尤度評価アルゴリズムは、いずれも式 (1), (2), および (3) に示すガウス分布の十分統計量を利用したものである。

$$A^0(m) = \sum_{t \in T} \gamma_m(t), \quad (1)$$

$$A^1(m) = \sum_{t \in T} \mathbf{x}_t \gamma_m(t), \quad (2)$$

$$A^2(m) = \sum_{t \in T} \mathbf{x}_t^2 \gamma_m(t). \quad (3)$$

ここで、 T は学習セット、 t は時刻 (フレームインデックス)、 m は混合要素のインデックス、 $\mathbf{x}_t =$

$(x_1(t), x_2(t), \dots, x_d(t))^T$ は時刻 t における d 次元特徴量ベクトル、 $\mathbf{x}^2 = (x_1^2, x_2^2, \dots, x_d^2)^T$ は特徴量の各要素を 2 乗したベクトル、 $\gamma_m(t) = P(m_t | T, \theta_0)$ は適当な初期モデル θ_0 に基づき計算された m 番目の要素分布の時刻 t における占有確率である。

混合ガウス分布 θ の自己尤度値は、最適化の過程においてアライメントが変化しないと仮定すると [1]、式 (5) のように表すことができる¹。

$$\begin{aligned} L_{self}(\theta) &\approx \sum_{m=1}^M \sum_{t \in T} \log(P(x_t | m, \theta)) \gamma_m(t) \quad (4) \\ &= -\frac{1}{2} \sum_m \left\{ \left(\log \left((2\pi)^d |\Sigma(m)| \right) + d \right) \cdot A^0(m) \right\}. \quad (5) \end{aligned}$$

ここで、 $\Sigma(m)$ は m 番目の要素ガウス分布の対角共分散行列である。この対角共分散行列は事前に計算された十分統計量を用いて式 (6) および (7) から求めることができる。

$$\boldsymbol{\mu}(m) = \frac{\mathbf{A}^1(m)}{A^0(m)}, \quad (6)$$

$$\begin{aligned} \text{diag}(\Sigma(m)) &= \mathbf{v}(m) \\ &= \frac{\mathbf{A}^2(m)}{A^0(m)} - \boldsymbol{\mu}(m)^2. \quad (7) \end{aligned}$$

さらに、併合操作後のガウス分布の十分統計量は併合前のガウス分布対の十分統計量を単純に足し合わせるにより容易に求められることから、最適化の全ての過程において、自己尤度を学習データを直接参照することなく事前に計算された十分統計量を用いて高速に計算することができる。

CV 尤度、および AgCV 尤度に基づく混合ガウス分布の最適化においては、まず学習セットを K 個の区画に区画化し、十分統計量を各区画ごとに計算しておく。ここで、 k 番目の区画に対して計算した十分統計量を $\mathbf{A}_k = \{A_k^0, \mathbf{A}_k^1, \mathbf{A}_k^2\}$ と表すことにする。

混合ガウス分布 θ の CV 尤度は自己尤度の計算と同じ仮定の下で、以下のように表すことができる。

$$L_{cv}(\theta) = \sum_{k=1}^K \sum_{m=1}^M \sum_{t \in T_k} \log(P(x_t | m, \theta_k)) \gamma_m(t). \quad (8)$$

ここで、 θ_k は k 番目のデータ区画を除いて学習した混合ガウス分布である。混合ガウス分布 θ_k のパラメタは k 番目の区画を除いて十分統計量を足し合わせることで、式 (9) および (10) に示すように容

¹混合分布重みはアライメント固定の仮定の下では混合分布最適化と独立であり、省略してある。

易に求めることができる。

$$\boldsymbol{\mu}_k(m) = \frac{\sum_{k \neq k} \mathbf{A}_k^1(m)}{\sum_{k \neq k} \mathbf{A}_k^0(m)}, \quad (9)$$

$$\mathbf{v}_k(m) = \frac{\sum_{k \neq k} \mathbf{A}_k^2(m)}{\sum_{k \neq k} \mathbf{A}_k^0(m)} - \boldsymbol{\mu}_k(m)^2. \quad (10)$$

式 (8) に、式 (9) および (10) により得られるモデルパラメタを代入し式変形することで、CV 尤度を事前に計算された十分統計量のみを用いて効率的に評価することができる [6]。

CV 尤度と同様に、AgCV 尤度は以下のように定式化できる。

$$L_{AgCV}(\theta) = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N \sum_{m=1}^M \sum_{t \in T_k} \{\log(P(x_t|m, \theta_{k,n})) \cdot \gamma_m(t)\}, \quad (11)$$

$$\boldsymbol{\mu}_{k,n}(m) = \frac{\sum_{i \in \Omega_{k,n}} \mathbf{A}_i^1(m)}{\sum_{i \in \Omega_{k,n}} \mathbf{A}_i^0(m)}, \quad (12)$$

$$\mathbf{v}_{k,n}(m) = \frac{\sum_{i \in \Omega_{k,n}} \mathbf{A}_i^2(m)}{\sum_{i \in \Omega_{k,n}} \mathbf{A}_i^0(m)} - \boldsymbol{\mu}_{k,n}(m)^2. \quad (13)$$

$\Omega_{k,n}$ は k を除いた 1 から K までの自然数 $(\{1, 2, \dots, K\} \setminus \{k\})$ の中から非復元抽出により無作為に K' 個の要素を選んだ集合である。ここで各 k について N 個の $\Omega_{k,n}$ の重複を許さない ($\Omega_{k,s} \neq \Omega_{k,t}$ if $s \neq t$) とすると、 N として用いることができる最大値は、 $C(i, j)$ を i 個のものの中から j 個の要素を取り出す組み合わせの数として、 $C(K-1, K')$ となる。式 (11) は CV 尤度法と同様な式変形により、学習データを直接参照することなく、事前に計算された十分統計量のみを用いて N に比例した計算量で効率的に計算することが可能である。

3.3 混合要素併合最適化の予備実験

図 3 に、ある HMM の状態に対し自己尤度法および AgCV 尤度法による最適化を適用した場合の、尤度変化の例を示す。横軸は混合数であり、初期状態の 200 から最適化の進行とともに減少していく。AgCV 尤度の計算では $K=6, N=10$ とした。図より、自己尤度はモデルの推定と尤度の評価に同じデータを用いているため正のバイアスが存在し、AgCV 尤度より大きな値をとるとともに、混合数に対して単調であることがわかる。これに対し AgCV 尤度はバイアスが大きく削減されているため、学習データに含まれない新しいデータに対するモデル性能のよい近

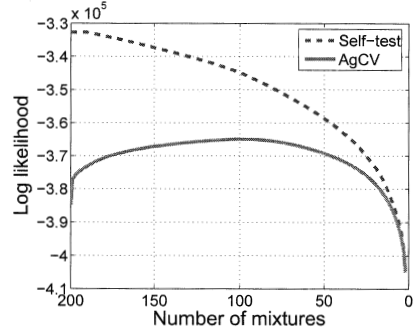


図 3: An example of the objective scores estimated for training data by the Gaussian mixture optimization methods.

似となり、混合数に対して最大値が存在する。この例の場合では、およそ 100 混合が最適であることが分かる。

4 実験条件

提案法を CSJ30 時間のデータを用いた混合ガウス分布 HMM の学習に適用した。評価セットは男性話者 10 名による学会講演から構成される CSJ のテストセットである。HMM の学習は 1 混合のトライフォンモデルを初期モデルとして入力し、EM5 回、提案法による各 HMM 状態の混合分布のモデル構造最適化、および全要素分布の二分割を合わせて 1 学習ループとして、15 学習ループ繰り返すことにより行った。固定された CV 区画に起因するバイアスを避けるため、各学習ループごとに学習ファイルリストの無作為化および学習区画の切り直しを行っている。詳細は文献 [6] と同様である。AgCV は $K=6, N=10$ とし、尤度最大を併合停止基準とした。

5 実験結果

図 4 に各学習ループにおいて 5 回の EM ループと提案法の AgCV によるモデル構造最適化に要した CPU 時間を示す。AgCV 法は CV 法と比較すると N に比例した計算量が必要となるが、実用的な計算量に収まっていることが分かる。

図 5 に学習ループを横軸にした単語誤り率を示す。図で “EM” はモデル構造最適化を行わないベースラインの結果である。“EM+MDL”, “EM+CVMIX”, “EM+AgCVMIX” はそれぞれ、MDL, CV, および本論文での提案法である AgCV を用いたモデ

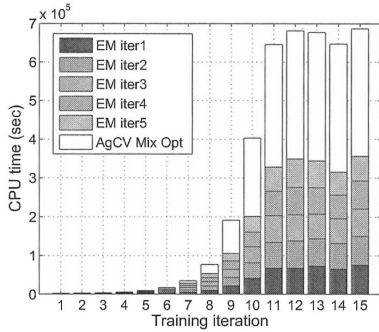


図 4: Computational cost.

ル構造最適化を行った場合の結果である。また、“AgEM+AgCVMIX”はEMの代わりにAgEM [8, 9]によるパラメタ推定と、AgCVによるモデル構造最適化を組み合わせた結果である。AgEMでは $K = 12, K' = 6, N = 12$ とした。

モデル構造の最適化を行わないベースラインの場合、単語誤り率は始めの間は学習ループ数とともに減少し、7回目の学習ループで最小値27.4%をとった後、再び増加に転じている。これはモデルのパラメタ数が学習ループ毎に単純に2倍に増えるため、学習ループがある程度以上に増えると、モデル推定精度の低下が顕著となるためである。

MDL, CV, および AgCV法を用いてモデル構造の最適化を行った場合、モデルのパラメタ数は自動的に制御されるため、単語誤り率は学習の繰り返しとともに次第にはば一定となる。また図より、AgCV法がMDL法やCV法による最適化よりも高い認識性能を与えることが分かる。これはAgCVがモデル選択法としてMDL法やCV法よりも優れていることを示している。また、EMの代わりにこれまでの提案法であるAgEMと本研究の提案法であるAgCVを組み合わせることで単語誤り率はさらに低下し、最低の誤り率26.2%が得られた。

6 まとめと課題

CV尤度による混合ガウス分布の最適化法を拡張した、AgCV法による混合ガウス分布の最適化法の提案を行い、大語彙連続音声認識実験により提案法がモデルサイズを自動決定し、更に従来法と比較して高い認識性能を与えることを示した。今後の課題としては、HMM状態クラスタリング等への応用が挙げられる。またAgCVはデータ駆動手法であり、十分統計量が利用できれば尤度以外の目的関数を用いたモデル学習へも応用可能と考えられる。

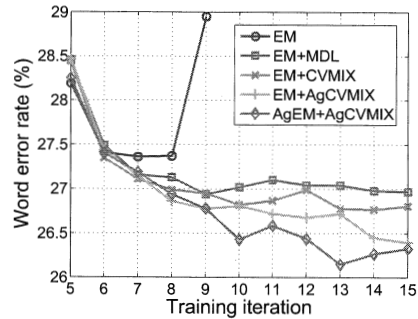


図 5: Number of training iterations and test set word error rates.

7 謝辞

本研究は科研費(19700167)の助成を受けたものである。

参考文献

- [1] S. Young, J. Odell, and P. Woodland. Tree-based state tying for high accuracy acoustic modelling. In *Proc. ARPA Workshop on Human Language Technology*, pages 307–312, 1994.
- [2] K. Shinoda and T. Watanabe. Acoustic modeling based on the MDL criterion for speech recognition. In *Proc. EuroSpeech*, volume 1, pages 99–102, 1997.
- [3] M. Ostendorf and H. Singer. HMM topology design using maximum likelihood successive state splitting. *Computer Speech and Language*, 11:17–41, 1997.
- [4] T. Cincarek, T. Tomoki, H. Saruwatari, and K. Shikano. Utterance-based selective training for the automatic creation of task-dependent acoustic models. *IEICE Transactions on Information and Systems*, E89-D(3):962–969, 2006.
- [5] T. Shinozaki. HMM state clustering based on efficient cross-validation. In *Proc. ICASSP*, volume I, pages 1157–1160, Toulouse, 2006.
- [6] T. Shinozaki and T. Kawahara. Gaussian mixture optimization for HMM based on efficient cross-validation. In *Proc. Interspeech*, pages 2061–2064, 2007.
- [7] L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
- [8] T. Shinozaki and T. Kawahara. GMM and HMM training by aggregated EM algorithm with increased ensemble sizes for robust parameter estimation. In *Proc. ICASSP*, pages 4405–4408, 2008.
- [9] T. Shinozaki and M. Ostendorf. Cross-validation and aggregated EM training for robust parameter estimation. *Computer speech and language*, 22(2):185–195, 2008.
- [10] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *J. of the Royal Statistical Society, Series B* 39(1):1–38, 1977.