

7. 音声の認識：高度な音響モデル

7.1 実際の音響モデル

7.2 識別的学習

7.3 深層学習

7.1 実際の音響モデル

- 混合分布の学習

- 各音素の特徴ベクトルは、一つの正規分布で近似できるほど単純ではない
例) 男女差、方言、...
- 複雑な確率密度関数を複数の正規分布の重み付き和で表現 → 混合分布

$$\phi = \sum_{i=1}^N w_i \phi_i$$

Φ_i : i 番目の正規分布
 w_i : i 番目の正規分布の重み
 N : 混合数

- 重みはEMアルゴリズムで学習

7.1 実際の音響モデル

- 話者適応

- 不特定話者用音響モデルのパラメータを、少数の特定話者データを用いて調整
- MLLR (Maximum Likelihood Linear Regression) 法
 - 学習済みHMMにおいて、平均ベクトルを以下の式で変換

$$\mu' = A\mu + b$$

- 特定話者データの尤度が最大となるような行列 A と定数項 b を推定

7.2 識別的学習

- 学習データの尤度計算

$$p(\mathbf{W}|\mathbf{X}) = \frac{P(\mathbf{X}|\mathbf{W})P(\mathbf{W})}{P(\mathbf{X})} = \frac{P(\mathbf{X}|\mathbf{W})P(\mathbf{W})}{\sum_{\mathbf{W}} P(\mathbf{X}|\mathbf{W})P(\mathbf{W})}$$

- 生成モデル： $P(\mathbf{X}|\mathbf{W})$ が大きくなるようにパラメータを求めた
- 識別モデルの考え方： $\sum_{\mathbf{W}} P(\mathbf{X}|\mathbf{W})P(\mathbf{W})$ を小さくすればよい
→ 正解以外の単語列に対して $P(\mathbf{X}|\mathbf{W})$ が小さくなるように学習

- 相互情報量最大化基準

$$\hat{\theta} = \arg \max_{\theta} \log P(\mathbf{W}|\mathbf{X})$$

$$= \arg \max_{\theta} \sum_r \log \frac{P(\mathbf{W}_r, \mathbf{X}_r; \theta)}{\sum_{\tilde{\mathbf{W}}} P(\tilde{\mathbf{W}}, \mathbf{X}_r; \theta)}$$

$\tilde{\mathbf{W}}$: 対立仮説
 r : 学習データの
インデックス

7.3 深層学習

- DNN-HMM法

- HMMの各状態で特徴ベクトルを出力する確率 $b_i(\mathbf{x})$ を $p(\mathbf{x} | s_i)$ と書き換え
- ベイズの定理

$$p(\mathbf{x} | s_i) = \frac{P(s_i | \mathbf{x})}{P(s_i)} p(\mathbf{x})$$

DNNで計算

学習データ
から最尤推定

定数

- \mathbf{x} はMFCCではなく、メルフィルタバンクの出力（またはもとの音声信号）で特徴抽出もDNNで学習

7.3 深層学習

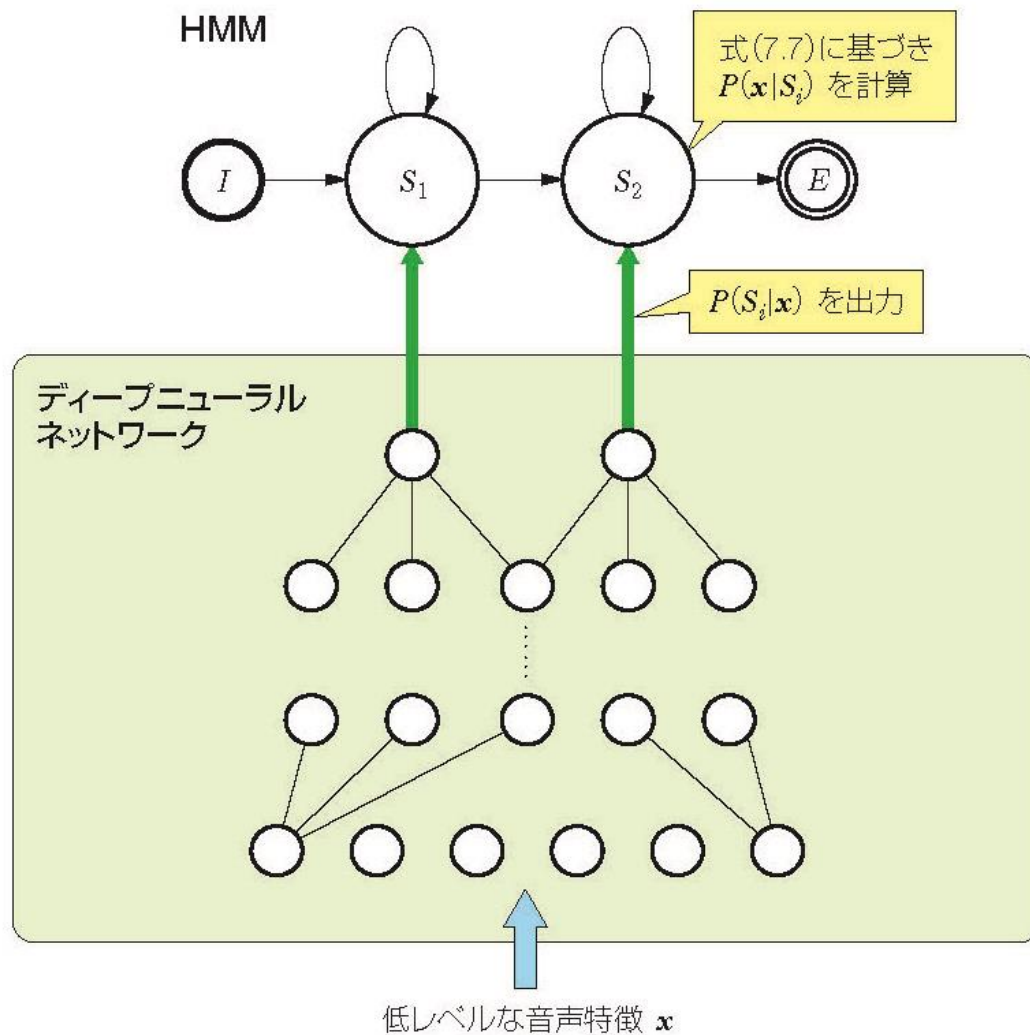


図 7.4 深層学習による音声認識