

第4章 統計的音声認識

4.1 枠組み

音声認識システムでは、音声による通信を、話者のメッセージ M が記号列 W に符号化され、それが音声波形 S に形を変えて通信路を経由し、聴き手に \tilde{S} として伝わり、 \tilde{S} の特徴分析で得られた観測時系列 O から記号列 \tilde{W} を復元し、その記号列から発話者のメッセージ \tilde{M} を推定する過程と捉えます (図 4.1).

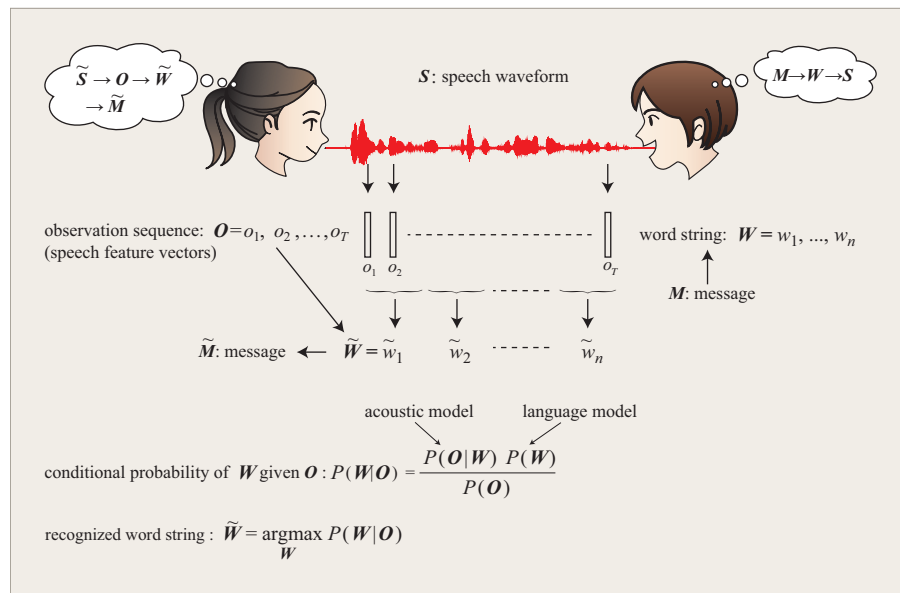


図 4.1: 音声通信の統計的解釈. M : メッセージ (話者が伝えたい内容), W : 単語列 (発話), S : 音声波形, \tilde{S} : 聴き手が受け取った音声波形, O : 観測系列 (聴き手の受け取った音声特徴の系列), \tilde{W} : 推定単語列, \tilde{M} : 推定メッセージ

音声認識の究極の目的は、音声波形 \tilde{S} から話者のメッセージ M の推定値 \tilde{M} を得ることであるともいえます。たとえば、部屋に人間とロボットが居るとして、人間が室温を下げて欲しいと思い、ロボットに「暑いですね。」と話したとします。その発話を聴いてロボットがエアコンの温度を調節して室温を下げてくれたならば、ロボットは人間のメッセージ (意図) を理解したということが出来ます。ロボットにこのような対応をさせるためには、単に聴いた音声波形を「暑いですね。」という文に変換することができるだけではダメです。発話の字面の背後にある意図を理解することが必要であり、音声理解 (speech understanding) と呼ばれています。現在、一般に音声認識と呼ばれている技術は、意図理解までは行わず、話者の発した単語列 W の推定値 \tilde{W} を求めることを指しています。

4.2 基本原理

最大事後確率基準に基づくベイズ (Bayes) の識別規則に基づいた音声認識の原理を説明します。

音響特徴量の系列を \mathbf{O} とします。 \mathbf{O} は、 m 次元の実数値ベクトル $\mathbf{o}_t \in R^m$ の長さ T の系列

$$\mathbf{O} \triangleq \mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T \quad (4.1)$$

とします。一般に、 \mathbf{O} は音声波形サンプル $\mathbf{X} = (x_1, \dots, x_N)$ に対して短時間フレーム分析を行うことにより得ます。発話された単語列を

$$\mathbf{W} \triangleq w_1, w_2, \dots, w_n, \quad w_i \in \mathcal{W} \quad (4.2)$$

とします。 \mathcal{W} は、ある有限の語彙 \mathcal{W} の要素 w_i からなる長さ n の単語列です。

$P(\mathbf{W}|\mathbf{O})$ は音響特徴量系列 \mathbf{O} が与えられたとき単語列 \mathbf{W} が発話された条件付き確率であり、事後確率 (posterior probability) と呼ばれます。ベイズの定理により、 $P(\mathbf{W}|\mathbf{O})$ は、

$$P(\mathbf{W}|\mathbf{O}) = \frac{P(\mathbf{O}|\mathbf{W})P(\mathbf{W})}{P(\mathbf{O})} \quad (4.3)$$

と書き直すことができます。

$P(\mathbf{W})$ は単語列 \mathbf{W} の生起確率です。たとえば、単語列 \mathbf{W} が 0.7 の確率で生起し、他の単語列をすべて合わせても 0.3 の確率でしか生起しなければ、どの観測データ \mathbf{O} に対してもそれが単語列 \mathbf{W} に属すると答えれば、0.7 の確率で正答となります。観測する前からわかっている確率という意味で、これを事前確率 (prior probability) と呼びます。

$P(\mathbf{O}|\mathbf{W})$ は単語列 \mathbf{W} が生起したという条件の下で音声特徴量時系列 \mathbf{O} が観測される確率を表わしていて、クラス条件付き確率 (class conditional probability) です。 \mathbf{O} が \mathbf{W} に属しているのが尤もらしいと考えられる確率と解釈することができることから尤度 (likelihood) とも呼ばれます¹。

$P(\mathbf{O})$ は \mathbf{O} が観測される確率であり、 \mathbf{O} と \mathbf{W} の同時確率 $P(\mathbf{O}, \mathbf{W})$ から

$$P(\mathbf{O}) = \sum_{\mathbf{W}} P(\mathbf{O}, \mathbf{W}) = \sum_{\mathbf{W}} P(\mathbf{O}|\mathbf{W})P(\mathbf{W}) \quad (4.4)$$

のように可能な単語列全体の関する和をとることで得られます。このような操作を周辺化といい、 $P(\mathbf{O})$ を周辺確率と呼びます。

式 4.3 を眺めてみると、事後確率 $P(\mathbf{W}|\mathbf{O})$ は事前確率 $P(\mathbf{W})$ をクラス条件付き確率と周辺確率の比 $P(\mathbf{O}|\mathbf{W})/P(\mathbf{O})$ で修正したものともみることができるとに気が付くでしょう。すなわち、単語列 \mathbf{W} が与えられたときに観測値 \mathbf{O} が得られる尤度 $P(\mathbf{O}|\mathbf{W})$ が、単語列 \mathbf{W} が与えられない場合の確率 $P(\mathbf{O})$ よりも大きければ、事後確率は事前確率よりも大きくなり、そうでなければ事後確率が事前確率よりも小さくなります。

¹ただし、尤度の和は $\sum_{\mathbf{W}} P(\mathbf{O}|\mathbf{W}) \neq 1$ なので、厳密な意味で確率とはいえません。

単語列 W_i と W_j の識別境界 (discrimination boundary) は, 事後確率が等しくなる場所, すなわち,

$$P(W_i|O) = \frac{P(O|W_i)P(W_i)}{P(O)} = \frac{P(O|W_j)P(W_j)}{P(O)} = P(W_j|O) \quad (4.5)$$

が成立するところです. 式 4.5 からわかるように, 周辺確率 $P(O)$ は単語列に共通に現れているので, 識別規則に含める必要はありません. したがって, 音声認識器は $P(O|W)P(W)$ を最大化する単語列 \tilde{W} ,

$$\tilde{W} \triangleq \underset{W}{\operatorname{argmax}} P(O|W)P(W) \quad (4.6)$$

を求めることになります. $P(O|W)$ を音響モデル (acoustic model), $P(W)$ を言語モデル (language model) と言います.