

文献 [3] では、音声認識は正しい認識結果を出力することを前提に (2) (3) を課題としている。しかし、後で示すように音声認識結果が88%程度でも分割精度やタグ付与精度はかなり劣化する。したがって、音声認識スコアだけを用いて、まず音声認識結果を求めてから、(2) (3) を実行したのでは精度のよい分割やタグ付与が期待できない。また、分割可能性確率やタグ付与確率の利用は音声認識精度そのものを向上させる可能性がある。つまり、課題は (1) (2) (3) を同時に満足する手法を発見することである。

3. 確率モデルと計算法

2章で提起した問題を解くために導入した確率モデルは基本的に文献[3]と同じものである。ただし、入力形態素列 W ではなく、候補単語ラティス V である。また、最適な形態素列 W (つまり音声認識結果)、最適な発話意図単位列 U 、最適なSAタグ列 T を同時に求めることになるので、式で書くと、

$$\begin{aligned} (\hat{U}, \hat{T}, \hat{W}) &= \arg \max_{U, T, W} P(U, T, W | V) \\ &= \arg \max_{U, T, W} \prod_{j=1}^k P(u_j, w_j, t_j | h_j, V) \end{aligned}$$

となる。ここで、 u_j はそのSA単位のj番目を示す。 h_j は u_j とj番目の t_j の履歴、すなわち (u_1^{j-1}, t_1^{j-1}) を示す。実際には1つ前、すなわち (u_{j-1}, t_{j-1}) だけで履歴を代表させている。SA単位の存在確率とタグ付与確率の計算法は文献[3]と同じものを採用している。

図1の各経路には音声認識スコア(音響スコアと言語スコア)、分割可能性確率(ただしスコア化)、SAタグの付与確率(ただしスコア化)が付与されている。これらは対数をとってあるのでスコアの総和が最大となる経路を見つければよいことになる。しかし、ラティスが大きいためそのままでは計算量が膨大となってしまう。そこで、音声認識スコアの大きいものから5つ(ビーム幅が

5)だけを計算対象とした。さらに計算時間の節減のため分割可能性確率が高いノードでだけ計算を実施することとした。探索には前向きDP+後ろ向きA*のアルゴリズム[4]を用いた。

4. 実験

実験に使用したのは974個の発声データであり、手作業により正解の形態素列(音声認識結果)の特定、SA単位への適性な分割、適正なSAタグの付与を実施した。正解の形態素列に対してSA単位への分割、SAタグ付与実験を実施したところ、SAタグ付与精度がかなり悪いことを発見した。原因を追求したところアルゴリズム(決定木など)の不備ではなく、手作業によるSAタグ付与のゆらぎであることが判明した。そこで、SAタグ付与基準の見直しを行い、元々30個あったタグの中で定義があいまいなもの、他との区別が容易でないものを除いて23個に縮小した。この結果、クロズドテストによる実験では、発話分割、SAタグ付与に関して100%の精度(適合率、再現率とも)が保たれることを確認した。

実験は用意した図1のようなワードグラフの入力データ974発声に対して、10回の交差確認テストを繰り返すことによって実施した。また、評価対象としたのは出力される1-bestのみである。

実験結果を表2に示す。手法の覽で「1:01」などと書いてあるのは、「音声認識スコア:分割タグ付与スコア(分割スコア+タグ付与スコア)」の重み(絶対値ではなく変位の大きさに対するもの)を示す。比較のために正解形態素列に対する音声認識精度、分割精度、SAタグ付与精度の値を表2の最上行に示す。この場合は音声認識スコアを使わずにまず音声認識結果を確定し、その後分割タグ付与スコアを用いて分割点とタグを決めたことに相当する。音声認識精度は当然100%である。分割精度は95%くらいで、SAタグ付与精度は85%程度である。つまり、このアルゴリズム

ムではこれが上限（それ以上の精度を期待できない）であることが分かる。なお、交差数とは正解の分割点をまたがる誤分割が行われる回数を示す。正解分割点の個数は約2000である。正解形態素列に対しての交差数は0であるから、分割ミスのすべては過剰な分割によることも分かる。なお、SAタグ付与精度は「分割もタグも正しい」割合を示したものであるため、必ず分割精度より低い値となる。

表2の「1best認識結果」は音声認識スコアで1位となったものに対して、後で分割タグ付与スコアを用いて分割点とタグを決めたものである。提案する手法（音声認識と分割・タグ付与を同時に行なう方法）はこれより良い結果を期待したものである。表2を見ると残念ながら期待したほどの改善はみられない。「1:0.01」付近でわずかに音声認識精度とタグ付与精度の向上が見られるが、統計的に有意なレベルではない。予想されることではあるが、分割タグ付与スコアの重みを増していくと音声認識精度が下がり、これに伴って分割精度やタグ付与精度も低下している。

5. おわりに

音声認識の結果として得られる単語ラティスに対する分割処理とタグ付与を同時に行ないつつ音

声認識結果を得る手法を提案した。単語ラティスを対象としても発話意図単位（SA単位）への分割、各SA単位へのSAタグ付与が可能であることを示すことができた。しかし、実験結果で示したように分割タグ付与精度の向上と音声認識精度の向上に関しては、期待した結果は得られなかった。しかし、まだ実験結果の検討は十分でない。例えばラティスのビーム幅を変えたらどうなるか、分割スコアとタグ付与スコアの寄与比率を変えればどうなるか、などの検討はこれからである。アルゴリズムの改良を進め、実験によって精度向上を図ることができれば改めて報告したい。

参考文献

- [1]M. Nagata and T. Morimoto, "An Information-theoretic model of discourse for next utterance type prediction" Transactions of Information Processing Society of Japan, Vol.35, No.6, pp.1050-1061(1994)
- [2]N. Reithinger and E. Maier, "Utilizing statistical dialogue act processing in verbmobil" Proceeding of the 33rd Annual Meeting of the ACL, pp.116-121(1995)
- [3]田中英輝、横尾昭男：確率モデルによる発話の最適分割と意図認識，言語処理学会第5回年次大会pp.259-262, (1999)
- [4]M. Nagata, "A stochastic Japanese morphological analyzer using a forward-DP and backward-A* N-best search algorithm" Proceedings of Coling94, pp.201-207(1994)

表2. 実験結果

手法	音声認識精度		分割精度		SAタグ付与精度		交差数
	再現率	適合率	再現率	適合率	再現率	適合率	
正解形態素	100.00	100.00	94.61	96.17	84.63	86.03	0
1best認識結果	87.76	87.63	61.66	59.02	57.13	54.69	162
1 : 0.001	87.78	87.66	61.66	58.99	57.19	54.72	155
1 : 0.01	87.84	87.70	61.66	58.92	57.25	54.71	145
1 : 0.1	87.29	87.34	61.31	59.18	56.79	54.82	147
1 : 1	83.23	83.44	49.19	44.49	45.59	41.24	169