CRF を用いたアクセント結合処理に対する誤り分析 とその改良に向けた考察*

印南圭祐,渡辺美知子,峯松信明(東大・新領域),広瀬啓吉(東大・情報理工)

1 はじめに

仮名漢字混じり文を入力として自然な音声を出力 する日本語テキスト音声合成を実現するためには,文 中のアクセント核位置を適切に推定する必要がある. 日本語の場合アクセント核位置は,孤立単語発声時と 文発声時とでは容易に変化する(アクセント結合). そのため,孤立発声時のアクセント核位置情報(辞書 に搭載されている語彙情報)と種々の言語情報より, 与えられた文に対して適切な位置情報を推定するモ ジュール開発が必須となる.従来,アクセント価や結 合様式などの属性を定義することで,規則としてアク セント結合を記述する方法 [1] が用いられてきた.こ れに対して筆者等は, CRF(条件付き確率場) を用い た統計的なアクセント結合処理手法を提案し,規則処 理よりも高い精度を実現することができた[2][3].こ の場合,学習データの増加や素性の再設計による精 度向上の可能性がある.本報では,先行研究の CRF を用いた処理モデルに関して,学習データ増加によ る推定精度の変化と,誤推定結果の分析,及び,それ らに基づく改良案について検討したので報告する.

2 CRF を用いた統計的アクセント結合処理の先行研究

筆者等の先行研究 [2][3] ではまず,単独ラベラを用いて,JNAS[4] 及び ATR 音素バランス 503 文のテキストに,アクセント句境界及び核の位置情報をラベリングさせた.これらの言語情報は,ラベラの方言 / アクセント感覚に依存することがあるため,単独ラベラに全てのラベリングを依頼した.そして,利用可能な読み上げ文 40 セット (4,108 文)を,学習用 / 評価用に 35 セット (3,581 文) / 5 セット (527 文) へと分割し,CRF++[5] を用いて,文中アクセント核位置を推定した.なお,核情報は形態素単位で推定した.用いた素性の詳細は [2][3] を参照して戴きたい.

結果を,規則に基づく推定精度と共に Table. 2 に示す.集計は形態素単位とアクセント句単位で行ない,後者では,主アクセントのみの正誤判定と,副次アクセントを含めた全てのアクセントに着目した判定を行なった.更にエラー解析として,下記の2種のアクセント句に着目し,これらの正答率も算出した.

- 単純なアクセント句 {自立語}+{付属語}から構成され,アクセント結合が1回だけ起こるアクセント句
- 名詞連続を含むアクセント句 2 語以上連続して名詞が出現するアクセント句

Table 1 学習データ量の変化に伴う誤推定の減少率 (1 セット当たり約 100 文)

17.00	ト数	形態素	すべての句	単純	名詞連続
69	I T 女X			十代	口的连领
10	20	18.5%	15.6%	9.8%	15.5%
20	30	21.0%	21.1%	15.2%	20.3%
30	40	9.5%	10.0%	7.7%	13.1%
40	50	10.5%	11.1%	0.0%	9.4%
50	60	10.8%	10.9%	25.0%	6.3%
60	65	3.9%	2.7%	14.8%	-2.2%

3 アクセントデータベースの拡張

3.1 CRF の学習データ量の増加

先行研究で作成したアクセントデータベースは,JNAS 読み上げテキストの 165 セット中 40 セットしか利用していなかった.そこで,同一ラベラに JNAS 文に同様のラベリングを行わせることで,データベースを拡張した.2008 年 1 月の時点で,70 セット (約7,200 文)が利用可能な状態である.これを元に学習データを 35 セット (約3,500 文)から 65 セット (約6,500 文)に増加,評価データは先行研究の評価実験と同一の文セットとして CRF によるアクセント核位置推定の実験を行なった.Table.2 最下段に結果を示す.これから分かるように,学習データ量の増加により,アクセント結合処理精度は全体的に向上した.

3.2 CRF の学習データ量の漸次的増加

学習量増加による推定精度改善の効果及びその飽 和性を見るために,学習データ量を漸次的に変化さ せて同様の実験を行なった.評価データは引き続き先 行研究の実験と同一の文セットとし, 学習データを 初期段階の 10 セット (約 1,000 文) から 10 セットず つ増加させていき,アクセント句単位の誤推定の減 少率を調べた.また,60セットから65セットまで増 やした場合の変化も集計した.なお,アクセント句 中の全てのアクセントの推定が適切な場合を正答と して評価した. 結果を Table. 1 に示す. 全体の傾向 として,学習データ量の増加に伴い,誤り減少率は 降下することが分かる(推定精度の飽和).特に60 65 の場合には単純な句以外の減少率は非常に低く なっている(逆に単純な句における推定精度が大きく 改善されている理由については現在検討中である). また,推定精度が伸び悩む名詞連続を含むアクセン ト句では,僅かではあるが誤りが増加する結果となっ た.以上の結果より,今後は学習データ量の増加より も,誤り分析に基づく素性の改良,及び,推定モデル の改良がより重要になると考えられる.

^{*}Error analysis of CRF-based accent sandhi estimation and some discussions for performance improvement, by K. Innami, M. Watanabe, N. Minematsu and K. Hirose (The University of Tokyo)

Table 2 各手法によるアクセント結合処理精度の比較

	形態素	すべての核			主核のみ		
	加速系	すべての句	単純	名詞連続	すべての句	単純	名詞連続
規則に基づく手法	%	76.4%	94.4%	73.5%	76.8%	94.5%	74.2%
CRF を用いた手法	95.5%	89.4%	95.5%	83.4%	91.7%	96.4%	85.6%
学習量を増加させた CRF	96.5%	91.9%	97.2%	86.6%	93.5%	97.7%	87.9%

Table 3 アクセント句の構成内容による誤推定数・正答率の比較

	すべての句	単純な句	名詞連続	数詞	カタカナ	副助詞	助詞・助動詞連続
該当アクセント句数	3533	822	688	421	347	128	430
誤答数	285	23	92	36	40	20	62
正答率	91.9%	97.2%	86.6%	91.4%	88.5%	84.4%	85.6%

4 CRF 出力ラベルの誤り分析

先行研究では実験結果の分析において,2種類のアクセント句に着目している.本稿では,誤り解析結果をモデル改良に反映させることを目的として,新たな構成のアクセント句に着目し推定精度を検証した.検討した各種の句の精度を Table.3 にまとめて示す.

4.1 数詞を含むアクセント句

数詞を含むアクセント句は,全体ではさほど推定精度が低くない(91.4%). しかし数詞は「名詞-数詞」という,名詞のサブカテゴリとして辞書登録されているため,Table. 4 のように名詞連続を含む句との重複が多く,そのような句の正答率は大幅に低下することになる(85.6%). また,数字表現を桁ごとに別の形態素として扱っているため,複数の数詞が連続して出現しやすい. よってこのような数詞が連続出現するアクセント句への対策が有効だと考える.

4.2 カタカナ表現を含むアクセント句

外来語などのカタカナ表現の語はその大半が名詞であるが, Table. 3 及び 4 から分かる通り, 名詞連続として出現した場合と, そうでない場合とで精度は殆ど変わらない. よって, 名詞接続時に生じる誤推定に対する解析よりも, カタカナ表現に起因する誤り解析を行なうべきであると考える.

4.3 副助詞を含むアクセント句

副助詞は出現する頻度こそ低いものの,これを含む句は正答率が非常に低い.更に,誤推定された285のアクセント句中に出現する副助詞の個数は形態素単位で21であるが,正しく処理されたものは3つしかなかった.今後,副助詞とその前後の表現について更に分析を行なう必要がある.

4.4 助詞・助動詞連続を含むアクセント句

単純な句で生じないような結合パターンとして,助詞または助動詞の出現の連続を考えたところ,副助詞を含む句に次いで低い推定精度であった.副助詞を含むアクセント句も含め,助詞・助動詞を含むアクセント結合処理に対する改善が必要と分かる.

Table 4 評価データ中の分析カテゴリの重複

	名詞連続	数詞&名詞	
該当句数	688	132	120
誤答数	92	19	14
正答率	86.6%	85.6%	88.3%

5 誤り分析に基づく改良の検討

以上のエラー分析から、名詞連続、及び、助詞・助動詞への対策が必要といえる。利用しているアクセントデータベース中の付属語は、自立語と違って孤立発声時のアクセント型の情報が登録されていない。そこでアクセント辞書の記述に従い、付属語に対して、孤立発声時のアクセント型を含む、アクセント結合に関する属性を追加することを検討している。具体的な属性としては、付属語が別の語と連接した際にアクセント核を持ちうるかどうかを示すような、アクセント結合後の結果に関する情報の付与を考えており、評価実験に向けてラベリング作業を行なっている。

6 まとめ

CRFによるアクセント結合処理は、学習データの増加によって改善させられるが、その向上には限界がある、そこでエラー分析の結果から、名詞連続を含む句と、助詞・助動詞等の付属語に対するアプローチが重要だと考える、今後は誤り分析の強化と、それに基づくデータベース・CRFの学習素性の改良を行ない、評価実験にて検証する予定である。

謝辞 本研究に全面的にご協力頂いた黒岩龍様と, 度々御助言を頂いた特定領域研究「日本語コーパス」 電子化辞書班の皆様に厚く感謝申し上げます.

参考文献

- [1] **匂坂, 佐藤, 信学論**, J66-D7, pp.849-856, 1983.
- [2] 黒岩他 ,信学技報 ,SP2006-174 ,pp.31-36 ,2007 .
- [3] 黒岩, "日本語音声合成のためのアクセント結合規則の改善とデータベースに基づく統計的アクセント処理", 東京大学大学院情報理工学系研究科電子情報学専攻修士論文 2007.
- [4] http://www.mibel.cs.tsukuba.ac.jp/jnas/
- [5] http://crfpp.sourceforge.net/