

日本文訂正支援システムにおける 6B-4 未知語訂正候補抽出方式

高木伸一郎 安田恒雄 島崎勝美 松岡浩司
NTT情報通信処理研究所

1. はじめに

日本文データベースや新聞記事データ等の作成においては、計算機に入力した多量の漢字かな混じり文に含まれる誤字等の誤りを検出、訂正する作業の大幅な省力化が求められている。このため、新聞記事の校正等の日本文訂正作業の省力化を目的として計算機による実用的なレベルでの誤り検出を行う日本文訂正支援システムREVISEを開発した。(参考文献1参照)本稿では、さらに高位機能として実現した検出誤りのレベルに応じた訂正候補抽出方式の概要、特に未知語により検出した誤りの訂正候補抽出方式について述べる。

2. 誤りの分類とREVISEにおける訂正候補抽出方式

新聞記事で検出された誤りのうち誤字、脱字、誤挿、送り仮名等の文脈に依存しない校正レベルの誤りについて分類すると、次のように大別される。(図1参照)

- (1) 表記の誤り(文字レベルの誤りや送り仮名、表記のゆらぎ)
- (2) 表現の誤り(複合語内承接誤りや同音異義語等)
- (3) 言語情報以外の一般常識に関する誤り(固有名詞誤り等)

総誤り数: 1454 (新聞記事13万文字)

文脈依存なし 55%		文脈依存 45%		
表記誤り 60%		表現誤り 30%	他	
スタイルチェック	文法チェック	承接関係チェック、自動注意喚起方式	同音異義検定等	知識ベース照合
30%	30%	20%	10%	10%
習慣性あり	習慣性なし			
(A) 言い換え表現抽出方式	(C) 未知語訂正候補抽出方式	(B) 注意語抽出方式		
→ 現REVISEの誤り検定方式による対象範囲				

図1 誤り分類と対応する訂正候補抽出方式

ここで、対応する誤り検定方式は、(1)では、習慣性のある誤りを検定するスタイルチェックと習慣性のない誤りを検定する文法チェックである。また(2)では、固有名詞・数詞に関する意味的な承接関係の有無を検定する承接関係チェックや表1に示すヒューリスティックなルールによるシステム自動注意喚起方式である。ここでは、複合語漢字列内の誤字のように文法的に検出されにくくかつ習慣性のない誤りに対応するため、複合語内の単語係り受け解析強化を行っている。このほか、同音異義語誤り検定、固有名詞誤り検定の開発を進めている。

The Correct Candidacy Extract Method on REVISE (Revision Support System for Japanese Verbal Error)
Shinichiro TAKAGI, Tsuneo YASUDA, Katsumi SHIMAZAKI, Kouji Matsuoka
NTT Communications and Information Processing Laboratory

表1 自動注意喚起方式および検定精度

項目	FA	FB	FC	FD	FE	計
検出数	26	0	0	124	126	276
ヒット数	10	0	0	22	32	64
ヒット率	38	0	0	18	25	23%

表記以外の誤り: 912中

FA: 係り受け無し1文字接辞 例: 参連(参画)
 FB: 係り受け無しひらがな接辞 例: 白げ(白さ)
 FC: 連続接辞(接頭辞+接尾辞) 例: 不荷(不可)
 FD: 連続接辞(接頭辞+接尾辞以外) 例: 講買層(講買層)
 FE: 係り受け無し1文字一般名詞・固有名詞 例: 日本経済(日本経済)

これらの誤り検定方式によって検出された誤りの特徴に応じた訂正処理を行うために、REVISEでは、次の3種類の訂正候補抽出方式を実現した。

(A) 言い換え表現抽出方式(言い換え表現導入過程: 図2参照)

習慣性のある表記誤りについては利用者辞書に予め誤り表記と対となる正解表記を登録し、誤り検定処理で使用する校正辞書にペアで組み込むことによって、誤り表記の単語が認定された場合に訂正候補を抽出する。現在の評価では、ほぼ100%の訂正候補抽出精度を実現している。

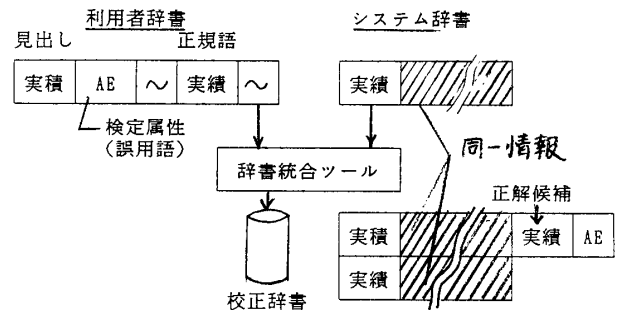


図2 言い換え表現訂正候補の導入法

(B) 注意語抽出方式

同音異義語や使用を誤り易い類義語・慣用語を注意語として予め注意語辞書に登録する事によりこれらの語が出現した場合には、利用者に対する注意喚起を行うとともに、注意語辞書の検索結果からの複数の訂正候補を抽出し提示する。

(C) 未知語訂正候補抽出方式

文法チェックにより検定される誤りは習慣性がなく、文法的に接続が不良な未知語を生成する特徴がある。このため未知語出現状況を分析して6種類の個別訂正候補抽出方式を考案し、その組み合わせによって有効な訂正処理機能の実現を図った。

3. 未知語訂正候補抽出方式

新聞記事1日分(約13万文字)に含まれる誤りで未知語によって検出された251件について未知語の出現特性を分類した次の6種類の個別訂正候補抽出方式を実現した。

(A) 高頻度類形文字・重複文字未知語訂正候補抽出

・高頻度で誤りやすい文字の対を予めテーブルに登録し、該当の文字が、未知語の前後に出現した際に置き換えて訂正候補を抽出する。

(例) ムルダ 二国軍 ----- ムルダニ (ニ---- ニ)

彼にとって ----- とって (つ---- っ)

・同一文字の誤挿により未知語化した際に重複文字の一方を削除する訂正候補を抽出する。

(例) 通じて て 得た ----- 通じて 得た

(B) 名詞+用言未知語訂正候補抽出

助詞抜け誤りにより見掛け上、「名詞+用言」化した場合に、文法的接続検定の結果発生する未知語を検出し、後方の用言が有する必須の格助詞を訂正候補として抽出する。必須格が複数存在する場合の訂正候補選択に関しては、前方の名詞句の意味カテゴリーについて照合する必要があるため、今後の機能としている。

(例) 決断 を 迫る ----- 決断 を 迫る (迫る--- を、に)

(C) 複合単語訂正候補抽出 (図3参照)

漢字列複合語は、一般に2文字単語の連続として構成される頻度が高いことを利用し、予め前後1文字のキーで検索できる高頻度漢字2文字候補テーブルを作成する。そして漢字列複合語内未知語の前後の1文字をキーとしてテーブルを索引し訂正候補を抽出する。この際、候補数が多数となるため、前後との文法的、意味的な接続条件を用いて候補の絞り込みを行う。

(例) 行 販 政 ----- 行 財 政 (財政)

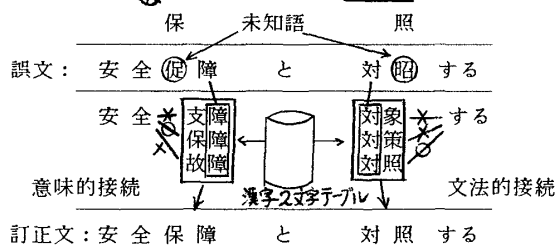


図3 複合語内漢字候補抽出方式

(D) 付属語尾文字訂正候補抽出

用言、用言性名詞での付属語尾が未知語となる場合に、語幹の活用型より付属語尾を抽出し訂正候補とする。さらに後方の単語の品詞により文法的に接続可能な活用形を選択する。

(例) 埋込 む ----- 埋め込 む (埋---- め)

(E) カタカナ・英字語訂正候補抽出

未知語を含むカタカナあるいは英字の文字列は見出し長も長く類似の単語も多いので誤りが混入しやすい。従ってカタカナあるいは英字辞書と照合をとることが有効であるが、対象単語数が多いので、2段階の照合過程による抽出方式を実現した。第1は、予め作成したカタカナ語訂正候補テーブルを2文字をキーに、少なくとも2文字以上マッチするカタカナあるいは英字語候補群を抽出する。つぎに高速に照合を行う連想統合照合法を用いて照合値を算出し脱字・誤挿補完処理を行って訂正候補順位を付与する。(図4参照)

本方式によって、48事例のカタカナ誤りに対し、正解候補含有率は1位85%、5位以内92%と非常に高い精度が得られた。

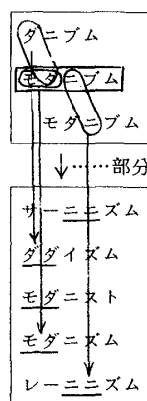


図4 カタカナ・英字訂正候補抽出方式

未知語: モダニズム (正解: モダニズム)
基本照合値テーブル

(F) 類音・類形文字(入力装置特性)訂正候補抽出

ペンタッチ・OCR等の入力特性に応じた誤りやすい類似文字を予めテーブル化し、未知語文字を用いて検索し候補抽出を行う。

(例) 諸国間 わ 関係 ----- の (わ---- の ペンタッチ)

懸案事項 ----- 懸案 (懸---- 懸 OCR)

4. 訂正候補抽出方式の評価

検討対象とした251件の未知語誤りに対する訂正候補抽出精度の評価結果を表2に示す。この結果、候補抽出処理の起動可能割合が82%、このうち候補数を考慮せず正解候補を抽出可能な割合は53%、さらに上位約10位以内に正解候補を含有する割合が49%となり、ほぼ未知語により検出した誤りの2件中1件は自動的に訂正候補が抽出できることがわかった。また起動されない誤りの6割はひらがな1文字誤字であることから、現在、文字の接続確率等の統計情報を用いる手法を検討している。

表2 未知語訂正候補抽出方式 評価結果

対象未知語数: 251件 (誤りを含む未知語件数) 約13万文字

訂正候補抽出方式	A	B	C	D	E	F	総計
候補抽出件数	20	4	47	22	92	21	206 (82.1%)
正解候補抽出件数	19	3	31	5	54	21	133 (53.0%)
上位候補選択件数*	19	3	22	5	52	11	122 (48.6%)

*訂正候補絞り込みの結果、正解が上位に残存するもの (上位約10候補以内)

5. おわりに

概述した訂正候補抽出方式の導入により、校正者は表記誤りについては計算機が示す候補の中から適当な候補を選択するというわずかの処理で校正できるので、省力化・文書作成の効率化を推進できる。今後は訂正精度を高めるとともに、同音異義語検定や知識ベース照合を用いた地名固有名称実在性検定等の誤り検定・訂正機能の検討を進める。

6. 参考文献

- 安田ほか:「日本文訂正支援システム REVISSE」
情処第33回全国大会 4J-9 1986,10
- 高木ほか:「日本文訂正支援システム REVISSE における誤り検定方式の検討」
情処第34回全国大会 6X-4 1987,3