

くずし字翻刻のための古文用音声認識システムの検討

研究代表者

文学研究科

情報理工学研究所

戸塚史織

ZHANG Yutao

1. 研究背景

くずし字

- 古典籍、古文書などの前近代資料に使われてきた文字
- 日本には数百万の古典籍・古文書があり、近年多くの機関がインターネットを通じてこれを公開し始めているが、くずし字で記されているため現代の多くの人は読めない

くずし字翻刻

- ・くずし字を現在で使用されている文字に変換する



舞鶴市糸井文庫所蔵
「花供養」の元資料

翻刻結果

- ・負荷が大きい
- ・くずし字解読では近年AIやOCR(光学的文字認識)技術が注目されている

現状

- くずし字解読を得意とする人々の多くはAIやOCR技術よりも、使い慣れたアナログ的な方法の方がより効率的に翻刻できる
- 翻刻に慣れた人が過去に自分の手元に手書きで残した翻刻のテキストデータ化が求められる

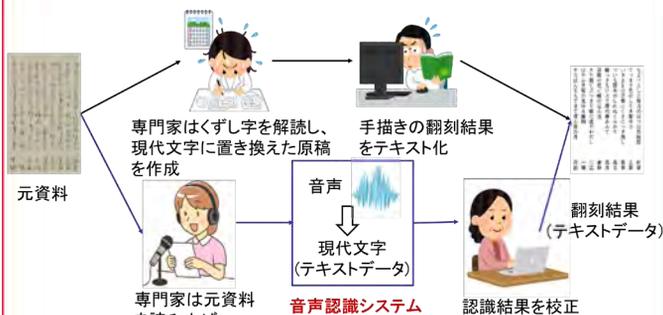
2. 古文用音声認識システム

くずし字翻刻の音声入力システム

- 負荷を減らし、翻刻のテキストデータ化効率化させるため
- ・画像だけではなく、音声からくずし字翻刻する方法を検討
- ・翻刻に音声認識システムの導入を検討

キープポイント:

高精度な古文認識を実現できる音声認識システムが求められる



古文用音声認識システムの構築

- 音声認識手法: DNN-HMM, End-to-End
- DNN-HMM古文用音声認識システム
 - ・認識エンジン: 大語彙連続音声認識エンジンJulius
 - ・音響モデル: ASJ-JNAS コーパスおよび CSJ(日本語話し言葉コーパス)で学習されたDNN-HMM音響モデルを利用

しかし、一般的な文章の言語データベースから学習された汎用言語モデルと発音辞書は古文に適用された場合、期待するほどの認識精度が実現できない



古文に対して認識精度を改善するため

古文用音声認識システムの言語モデルと発音辞書を構築

学習データの準備

- ・『日本語歴史コーパス』に含まれる作品
 - ◆ 江戸時代洒落本(15作品)、文字数: 約17万
- ・独自の学習データ
 - ◆ 役者白虎通 江戸の巻(先頭5000語)
 - ◆ 電子辞書の解析結果を元に修正

聖遊廊	新月花余情	花街弄々女(前編余興花街弄々女)	花街籠(玉菊全伝花街弄々女)	吳那中奇譚
箱まくら(河東方言箱まくら)	色深猿睡夢	風流様人形	深川新話	原柳巷花話
当世左様候	郭中奇譚	無論里問答	陽台道頓・嵯峨秘言	月花余情

言語モデルの構築

- ・形態素解析結果による
 - ◆ 学習データの処理: 文に分割する、不要部分を削除するなど
 - ◆ 辞書の構築: 語彙へ読みを付与する
- ・処理した学習データからn-gramを推定する
- ・バックオフn-gramモデルを構築する

3. 評価実験

実験目的

- 現時点で構築した言語モデルの性能を確認

実験手順

- テストデータ
 - ・洒落本の読み上げデータ 内容: 「聖遊廊」の一部 長さ: 4分くらい
- 音声認識ソフトウェアJuliusを用い、次の三つの言語モデルとその中に含まれた単語で構成された発音辞書でテストデータを認識
 - ・汎用言語モデル(『現代日本語書き言葉均衡コーパス』で学習した言語モデル)
 - ・クローズドテスト(テストデータ「聖遊廊」を含んでいる学習データで学習した言語モデル)
 - ◆ モデルが学習できたかどうかを確認
 - ・オープンテスト(テストデータ「聖遊廊」以外の学習データで学習した言語モデル)
 - ◆ 学習データに含まれない古文に対し、認識性能を確認

各言語モデルの文字誤り率



認識結果が良い例と悪い例

原文	心をやわらぐるもゆかりのつきの一ふしぞかし
汎用言語モデル	心を、和らぐ置ゆかりのお付きの人藤井浴かし
クローズドテスト	心をやわらぐるもゆかりのつきの一ふしぞかし
オープンテスト	心を和らぐるもゆかりの月のひとふしぞかし

原文	もしにさま此中横断でお見上げ申ましたゆへ
汎用言語モデル	あつせんさま、この、魔界より、中井を公算に寮を三重県申しましたゆえ
クローズドテスト	ます天さまあ此あつて間かみよりあつて中居を古語にてお見上げ申ましたゆへ
オープンテスト	ます天さまあ此あつて間かみよりあつて中居をこりてを身請の申ましたゆへ