

Confusion Network Based Vector Space Model を 用いた音声ドキュメントの検索*

○胡 新輝, 呉 友政, 柏岡 秀紀 (NiCT/ATR)

1 はじめに

本論文では、Confusion Network (CN) をベースにしたベクトル空間モデルを用いて、音声ドキュメントの検索手法を提案する。Confidence に依存するtfの計算と出現頻度に依存するidfの計算を結合する事によって、良い検索性能を得られる事が我々の研究によって確認された[1]。本論文では、これらの計算方法を使って、異なる音声認識率の結果に対して検索実験を行う。日本語話し言葉コーパス (CSJ) を対象にした実験で、本手法は、高い認識誤り率 (WER) の場合でも、1-bestの認識結果を用いる場合と比べて、6%以上の性能改善が得られた。

2 Confusion Network をベースにした ベクトル空間モデル

2.1 CN の導入

近年、目覚ましく進歩した音声認識技術のおかげで、1-best の結果でも、ニュース放送など比較的整った音声の検索は、満足度の高い検索性能が得られる[2]。しかし、自由発話の音声では、認識性能がなかなか改善できなく、単語の誤り率 (WER) が30%より大きい場合が多い。このような場合、失われた情報が多いため、1-best のみで検索を行うと、十分な検索性能が得られない。

この問題に対処するために、多数の認識の仮説を活用することを考える。音声認識のLatticeを使えば、Lattice中に含まれる大量の不確定な内容を利用することができ、1-bestで失われた有用な情報を補うことができると考えられる。我々はMangu [3] が提案したCNを採用して、自由発話の音声データに対して、情報検索の研究を試みている。

音声認識器が出力する認識結果のLatticeに変換をかけて[2], 図1のような単語CNを得る。このようにCNは、多数のalignment (図のノードに相当) と、alignment間に複数の競

合する単語候補を含んでいる。これらの候補はconfusion setという。Confusion setの各単語は、それぞれの事後確率を持つ。以下、CNに対してベクトル空間モデルを定義する。

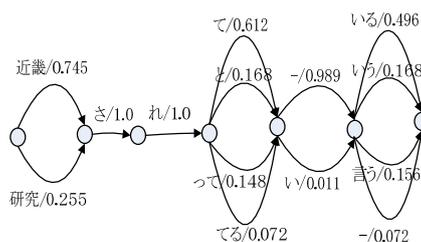


Fig. 1 CN の例

2.2 ベクトル空間モデルの定義

(1) 基本定義

D — CNが表現するドキュメント。

$P(w|o, D)$ — ドキュメントDの位置oに単語wが発生する事後確率。

(2) tf・idfのベクトル空間モデルの計算

我々は、テキスト検索によく使われているベクトル空間モデルであるtf・idf方式をCNに適用して、音声ドキュメントのキーワードの検索を行った。実験の結果として、以下の結論を得た[1]。

- tfの計算の際、confidenceが単語のランクより有効である。
 - idf計算の際、単純な単語カウントの方がconfidenceの使用より有効である。
- 従って、以下の様に、tfとidfを計算する。

$$tf(w, D) = \sum_{i=1}^{|\text{occ}(w, D)|} P(w | o_i, D) \quad (1)$$

$$idf(w) = \log(N / \sum_{D \in C} O(w, D)) \quad (2)$$

ここで、 $o(w, D) = \begin{cases} 1 & \text{if } P(w|o, D) > 0 \\ 0 & \text{otherwise} \end{cases}$

Cは、すべてのドキュメントの集合である。

(3) 検索 Query とドキュメント間の類似度の計算は、下式で行う。

* Spoken Document Retrieval Using Vector Space Model Based on Confusion Network, Xinhui Hu, Youzheng Wu, Hideki Kashioka (NiCT/ATR).

$$\text{sim}(d_i, q) = \frac{\sum_{w_q \in V} \text{tf}(w_q | d_i) \cdot \text{C}(w | q)}{\sqrt{0.8 \cdot \text{avdl} + 0.2 \cdot |d|}} \quad (3)$$

ここで、 $c(w|q)$ は、検索クエリ (Query) に現われた単語 w の頻度である。 $|d|$ は、ドキュメント d に含まれる CN の数である。

$\text{avdl} \left(= \frac{1}{|C|} \sum_{d \in C} |d| \right)$ はドキュメント長の平均である。

3 実験

3.1 データと実験条件

日本語話し言葉コーパス (CSJ) の「学会講演」と「模擬講演」の 2702 音声ファイルを使って、検索実験を行った。両方の講演を合わせて 600 時間を超える。上述の音声データに二種類の認識条件で認識を行い、それぞれ単語誤り率 (WER) が $\text{wer}_h=42.13$ と $\text{wer}_l=30.25$ の認識結果の CN を得た。

3.2 評価尺度

検索性能の評価尺度には、0.0 から 1.0 まで 0.1 刻みの各再現率における精度、それらの 11 点の再現率レベルにおける補間精度 IP を平均した 11 点の平均精度 AP、及び AP を全検索クエリで平均した MAP (Mean Average Precision) を用いた。

3.3 Query と適合性判定

Query 文は、CSJ テストコレクション [4] の中から解答が存在する 39 個の質問文に対して、形態素解析をした後、人手でストップワード、名詞以外のものを排除して作った。また、今回の検索は、既知語の検索であるため、Query のキーワード列に、未知語が含まれないようにした。適合性の判定作業は、TREC の評価ツールの `trec_eval` で行った。検索対象の単位としては、講演単位である。実際の判定は、適合 (Relevant) と不適合 (Irrelevant) の二種類で行い、テストコレクションの“部分適合”は適合と見なした。

3.4 実験結果

図2には、二つ認識率の条件での検索結果を示す。それぞれ 1-best に対する検索結果と比較している。どちらの認識率の場合でも、CN の使用により、検索性能が改善している。また、 wer_l と wer_h のそれぞれの改善は 2.2% と 6.2% であり、認識誤りが高い時の改善が大きい。

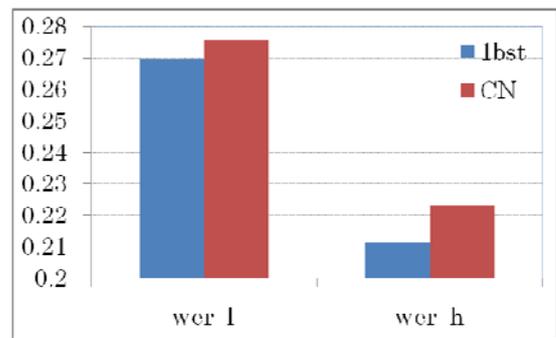


Fig. 2 異なる認識率時の検索性能

4 まとめ

本研究では、テキスト検索によく使われているベクトル空間モデルである $\text{tf} \cdot \text{idf}$ 方式を CN に適用する方法を提案し、それを用いて異なる音声認識率の音声ドキュメントの検索を作った。実験の結果として、以下の結論を得た。

- CN をベースにしたベクトル空間モデルは、1-best の結果より、検索性能が良い。
- 認識率が悪い時、CN は効果を発揮する。今後、より多くの認識率の異なる音声データを使って、実験を行う予定である。また、単語以外の検索単位も考慮して、認識誤りにより全体の性能の悪化に対する対策を研究する。

謝辞

本研究は、文部科学省科研補助金「特定領域研究」の“情報爆発時代に向けた新しい IT 基盤技術の研究” — “音声ドキュメントの検索に関する研究”により実施したものである。

参考文献

- [1] 胡, et al. “Lattice をベースにした音声ドキュメントの検索に関する研究”, 音講論集, 1-Q-32, Mar. 2008.
- [2] Garofolo, et al. “The TREC Spoken Document Retrieval Track: A Success Story.” Proc. RIAO Conference 2000: pp 1-20, 2000.
- [3] L. Mangu, et al., “Finding Consensus in Speech Recognition: Word Error Minimization and Other Applications of Confusion Networks,” Computer Speech and Language, vol. 14, no.4, pp.373-400, Oct, 2000.
- [4] 秋葉友良, et al. “音声ドキュメント検索テストコレクションの試作と検索性能評価,” 第1回音声ドキュメント処理ワークショップ, 2007.