

PT4-17 Masked Language Model による系列確率に基づく文法誤り検出

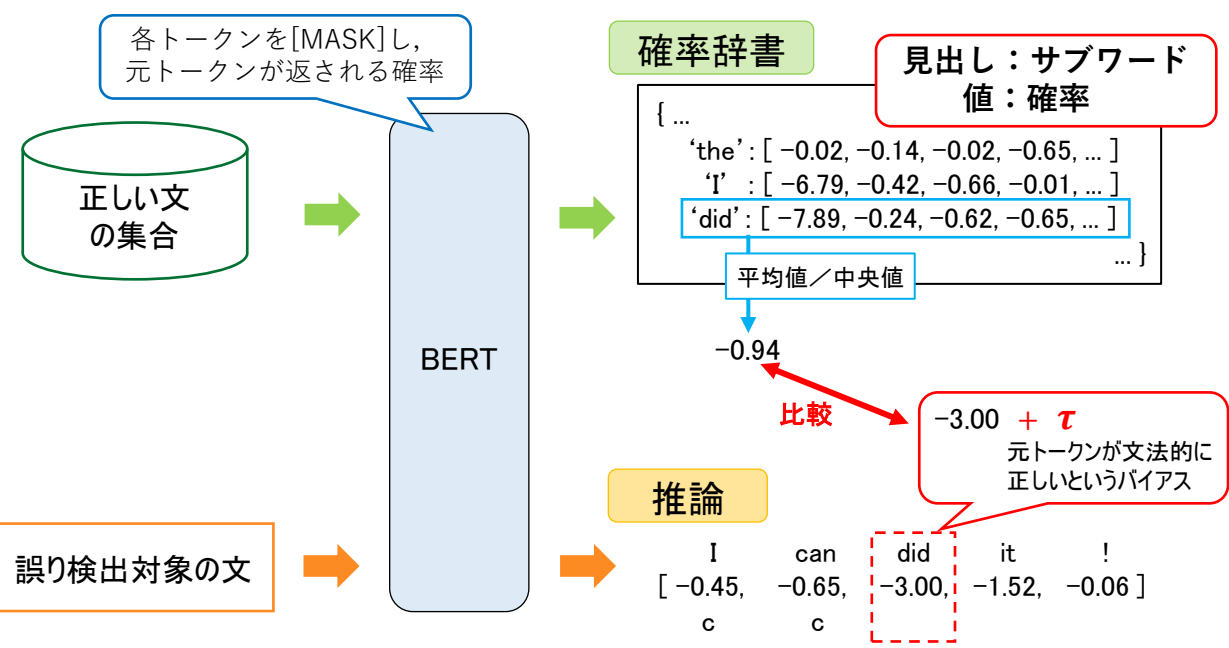
○土肥 康輔, 須藤 克仁, 中村 哲 (NAIST)

本研究の概要

BERTの確率を利用した**文法誤り検出**
 → **ラベル付きデータなし**で, **ナイーブなLSTM**をやや上回った

提案手法

- 確率の低い語は, 誤っているだろう
 ある文での“did”の確率を, “did”全体の確率と比較 → **[確率辞書]**

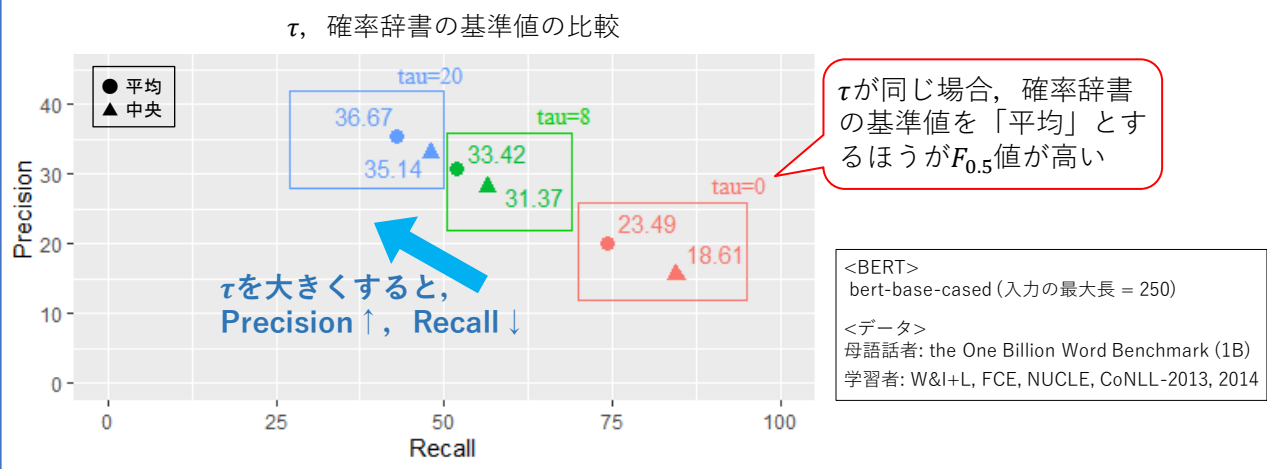


• **サブワードと元トークンの対応付け**
 すべて“c”のときのみ“c”

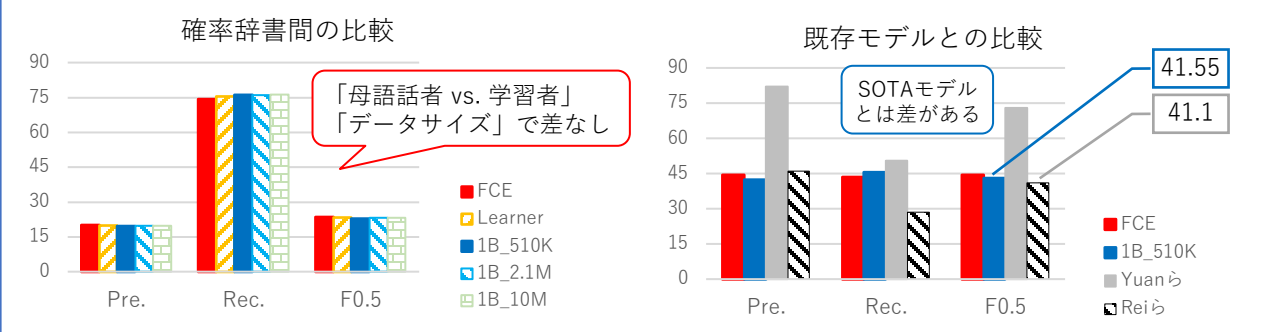
le	##mma	##tized	lemmatized
[c]	[c]	[c]	[c]
[c]	[i]	[c]	[i]

実験

確率辞書作成方法, τ の違いによって, 誤り検出性能に差があるか?



→ Recallが高くなる「中央」からスタートし, 適切な制約を加えるほうがよい



→ 母語話者コーパスで◎

今後の課題

- 文脈情報を利用した制約
- 文脈, 品詞等を考慮した確率辞書作成

- Precisionに課題
- ラベル付きデータなしで, **ナイーブなLSTM**をやや上回った

【参考文献】
 [Yuan+ 2021] Multi-class grammatical error detection for correction: A tale of two systems. In *Proc. of EMNLP*.
 [Rei+ 2016] Compositional sequence labeling models for error detection in learner writing. In *Proc. of ACL*.