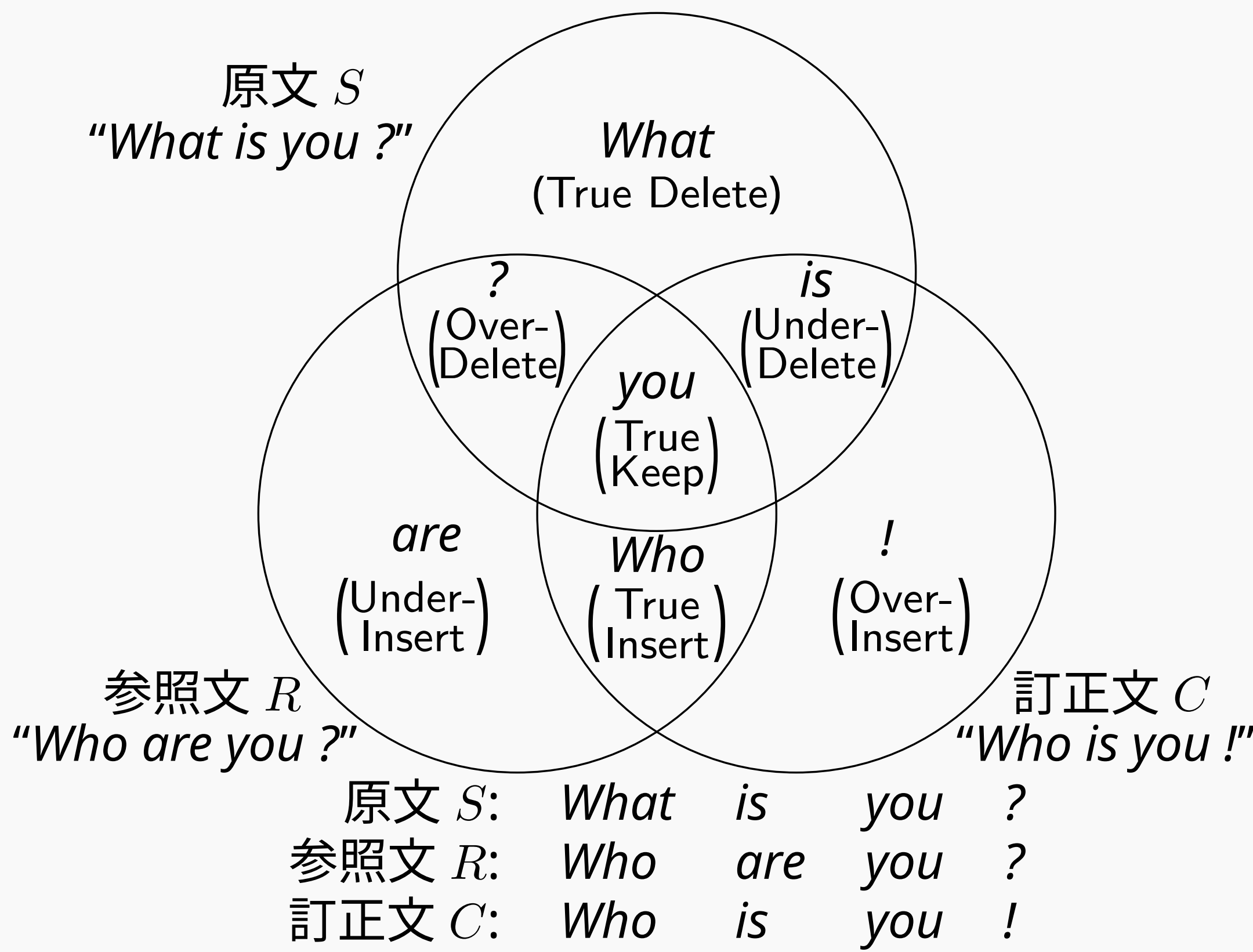


概要

- ▶ 文法誤り訂正とは？文書中の誤りを訂正するタスク。
 - ▶ 例: “What is you?” → “Who are you?”
- ▶ システムの出力への良い自動評価手法がほしい。
 - ▶ システムの質を手で評価するのは避けたい。
- ▶ M^2 とその派生が広く使われている。
 - ▶ 文間のフレーズアライメントを取り、訂正の個数を数え、 F -score を計算。
- ▶ M^2 系の問題点は？
 - ▶ アライメントの計算に時間がかかる。また、アノテーションされたデータや、その他のツールが必要となる。
- ▶ 提案手法 GREEN
 - ▶ 文を n -gram の多重集合とみなし、文間の編集を集合操作で抽出する。
 - ▶ アライメントの計算が不要。
 - ▶ 原文から参照文の編集、原文から訂正文の編集を比較し、 F -score を計算。
 - ▶ 既存の手法より人手評価と高い相関を示した。

手法

- ▶ 提案手法 GREEN は、文を最大長 N までの n -gram の多重集合として扱う。(例: 文 $A = “a a b”$ は、 $N = 2$ で、 $\{a, a, b, a a, a b\}$ 。)
- ▶ 原文 S , 参照文 R , 訂正文 C を多重集合で表し、ベン図を用いて文中の n -gram に対して 7 分類を行う。



領域	名前	$S \rightarrow R$	$S \rightarrow C$
$S \cap \bar{R} \cap \bar{C}$	True Delete (TD)	Delete	Delete
$\bar{S} \cap R \cap C$	True Insert (TI)	Insert	Insert
$S \cap R \cap C$	True Keep (TK)	Keep	Keep
$S \cap R \cap \bar{C}$	Over-Delete (OD)	Keep	Delete
$\bar{S} \cap \bar{R} \cap C$	Over-Insert (OI)	なし	Insert
$S \cap \bar{R} \cap C$	Under-Delete (UD)	Delete	Keep
$\bar{S} \cap R \cap \bar{C}$	Under-Insert (UI)	Insert	なし

- ▶ 多重集合の重複度（頻度） $m_A(x)$ を用いる。(例: $m_A(a) = 2$)
- ▶ 多重集合上の演算（積 \cap 、和 \cup 、差 \setminus ）を定義。

$$m_{A \cap B}(x) = \min(m_A(x), m_B(x))$$
$$m_{A \cup B}(x) = \max(m_A(x), m_B(x))$$
$$m_{A \setminus B}(x) = \max(m_A(x) - m_B(x), 0)$$

- ▶ 原文 S , 参照文 R , 訂正文 C の各 n -gram x を、ベン図を用いた 7 分類で数え上げる。

$$TD_{S,R,C}(x) = m_{S \cap \bar{R} \cap \bar{C}}(x) = m_{S \setminus (R \cup C)}(x)$$
$$= \max\{m_S(x) - \max(m_R(x), m_C(x)), 0\}$$

$$TI_{S,R,C}(x) = m_{\bar{S} \cap R \cap C}(x) = m_{(R \cap C) \setminus S}(x)$$
$$= \max\{\min(m_R(x), m_C(x)) - m_S(x), 0\}$$

$$TK_{S,R,C}(x) = m_{S \cap R \cap C}(x) = \min(m_S(x), m_R(x), m_C(x))$$

$$OD_{S,R,C}(x) = m_{S \cap R \cap \bar{C}}(x) = m_{(S \cap R) \setminus C}(x)$$
$$= \max\{\min(m_S(x), m_R(x)) - m_C(x), 0\}$$

$$OI_{S,R,C}(x) = m_{\bar{S} \cap \bar{R} \cap C}(x) = m_{C \setminus (S \cup R)}(x)$$
$$= \max\{m_C(x) - \max(m_S(x), m_R(x)), 0\}$$

$$UD_{S,R,C}(x) = m_{S \cap \bar{R} \cap C}(x) = m_{(S \cap C) \setminus R}(x)$$
$$= \max\{\min(m_S(x), m_C(x)) - m_R(x), 0\}$$

$$UI_{S,R,C}(x) = m_{\bar{S} \cap R \cap \bar{C}}(x) = m_{R \setminus (S \cup C)}(x)$$
$$= \max\{m_R(x) - \max(m_S(x), m_C(x)), 0\}$$

- ▶ 各 n で S, R, C ごとの True Positive (TP), False Positive (FP), False Negative (FN) を計算する。

$$TP_{n,S,R,C} = \sum_{\forall n\text{-gram } x} (TD_{S,R,C}(x) + TI_{S,R,C}(x) + TK_{S,R,C}(x))$$

$$FP_{n,S,R,C} = \sum_{\forall n\text{-gram } x} (OD_{S,R,C}(x) + OI_{S,R,C}(x))$$

$$FN_{n,S,R,C} = \sum_{\forall n\text{-gram } x} (UD_{S,R,C}(x) + UI_{S,R,C}(x))$$

- ▶ 各 n でコーパス全体 ($\mathcal{S} = (S_1, \dots), \mathcal{R} = (R_1, \dots), \mathcal{C} = (C_1, \dots)$) での precision, recall の幾何平均を取り、GREEN の F_β 値を求める。

$$\text{precision}(N, \mathcal{S}, \mathcal{R}, \mathcal{C}) = \left(\prod_{n=1}^N \frac{\sum_i TP_{n,S_i,R_i,C_i}}{\sum_i (TP_{n,S_i,R_i,C_i} + FP_{n,S_i,R_i,C_i})} \right)^{\frac{1}{N}}$$

$$\text{recall}(N, \mathcal{S}, \mathcal{R}, \mathcal{C}) = \left(\prod_{n=1}^N \frac{\sum_i TP_{n,S_i,R_i,C_i}}{\sum_i (TP_{n,S_i,R_i,C_i} + FN_{n,S_i,R_i,C_i})} \right)^{\frac{1}{N}}$$

$$\text{GREEN}_\beta(N, \mathcal{S}, \mathcal{R}, \mathcal{C}) = \frac{(1 + \beta^2) \text{precision}(N, \mathcal{S}, \mathcal{R}, \mathcal{C}) \text{recall}(N, \mathcal{S}, \mathcal{R}, \mathcal{C})}{\beta^2 \text{precision}(N, \mathcal{S}, \mathcal{R}, \mathcal{C}) + \text{recall}(N, \mathcal{S}, \mathcal{R}, \mathcal{C})}$$

- ▶ 複数の参照文 (R_{i_1}, \dots, R_{i_m}) がある場合は、文単位でスコアを最大にする \hat{R}_i を選ぶ。

$$\hat{R}_i = \operatorname{argmax}_{R \in \{R_{i_1}, \dots, R_{i_m}\}} \text{GREEN}_\beta(N, (S_i), (R), (C_i))$$

実験

- ▶ CoNLL-2014 評価タスクでの評価
 - ▶ Corpus-level 評価を行い、Expected Wins 法の**人手評価値との相関係数**、1 システムに対する**1 回の評価の実行時間**の平均を計測した。
 - ▶ アライメントによる F -score は、precision-oriented な評価を行うため、 $\beta = 0.5$ を用いる。 n -gram による手法は、参照データ数のバイアスを緩和するため、recall-oriented な評価を行い、 $\beta = 2.0$ を用いる。
 - ▶ n -gram の計測の分割は、単語単位、文字単位で行い、それぞれ、word+ 手法名、char+ 手法名と表記する。
 - ▶ charGREEN は、既存手法より高い相関を示す。
 - ▶ GREEN は、既存手法よりも高速に評価を行える。
 - ▶ M^2 は人手でアノテーションしたアライメントを必要とする。ERRANT、CLEME は品詞タグ付けによるアライメントを行う。PT- M^2 は BERT を用いる。GREEN は追加のデータやツールが不要かつ高速で、高性能である。

	r	ρ	時間 (秒)
アライメントによる F -score			
M^2	0.623	0.687	4.34
ERRANT	0.644	0.687	12.35
PT- M^2	0.686	0.786	109.82
CLEME	0.648	0.709	10.15
n -gram による precision			
wordGLEU	0.696	0.445	2.69
charGLEU	0.606	0.593	5.00
n -gram による F -score			
wordGREEN	0.741	0.698	0.55
charGREEN	0.786	0.813	1.94

- ▶ 評価手法ごとの違いは？
 - ▶ AMU : CoNLL-2014 タスクで最も人手評価値が高い出力
 - ▶ INPUT : CoNLL-2014 タスクの入力文をそのまま出力とする
 - ▶ IPN : CoNLL-2014 タスクで最も人手評価値が低い出力
 - ▶ NULL : すべての入力に対して、空の文を出力とする
- ▶ アライメントによる F -score は、悪い出力を高く評価してしまうことがある。

	AMU	INPUT	IPN	NULL
アライメントによる F -score				
M^2	35.01	0.00	7.09	28.01
ERRANT	31.97	0.00	5.95	0.20
PT- M^2	35.94	0.00	5.72	2.44
CLEME	25.14	0.00	4.41	33.44
n -gram による precision				
wordGLEU	58.08	56.34	55.08	0.00
n -gram による F -score				
wordGREEN	79.26	76.93	76.31	43.46
人手評価	62.84	45.64	29.99	-