An Empirical Study of Incorporating Pseudo Data into Grammatical Error Correction

Shun Kiyono^{1,2} Jun Suzuki^{2,1} Masato Mita^{1,2} Tomoya Mizumoto^{1,2*} Kentaro Inui^{2,1}
¹ RIKEN Center for Advanced Intelligence Project ² Tohoku University
{shun.kiyono, masato.mita, tomoya.mizumoto}@riken.jp;
{jun.suzuki,inui}@ecei.tohoku.ac.jp



Code available at https://github.com/butsugiri/gec-pseudodata

Grammatical Error Correction (GEC)

- Input: sentence with grammatical error
- Output: sentence without grammatical error
- GEC is commonly tackled as Machine Translation



Model (e.g. Encoder-Decoder)

GEC data is "low-resource"

- Amount of parallel data is limited in GEC (~2M)
- More data is important for better performance



Figure from [Sennrich and Zhang 2019]

- **Pseudo data generation** is currently the best method to increase GEC data
 - Adopted by most teams in BEA 2019 Shared Task

Training with Genuine Data only





Model





November 18, 2019

RIKEN AIP / Tohoku University

Problem: Lack of Consensus

• We have **three questions** regarding the pseudo data:

Q1: Choice for seed corpus? Q2: Methods for generating pseudo data? Q3: How to use pseudo data during training?

- GEC research community lacks consensus 谷
- Our aim: Find **settings that consistently improve performance** of GEC model



Q1: Choice for Seed Corpus?

- Numerous options exist:
 - Wikipedia, 1-billion word benchmark, BookCorpus, etc
 - [Ge+2018] uses Wikipedia
 - [Zhao+2019] uses 1-billion word LM benchmark
 - [Xie+2018] uses NYT corpus
 - [Grundkiewicz+2019] uses News Crawl
- What kind of corpus is suitable for GEC?
- We compare following three corpora:
 - Simple Wikipedia Grammatical complexity is different
 - Wikipedia
 - LDC Gigaword ——— Texts are cleaner in Gigaword

Q2: Methods for Generating Pseudo Data?

"Genuine" Data



Q2: Methods for Generating Pseudo Data?

Original:	At the in00 stitute , she introduced tis00 sue culture
	methods that she had learned in the U.00 S.
BACKTRANS (NOISY)	:At in@@ stitute , She introduced tis@@ sue culture method
DIRECTNOISE:	that she learned in U.00 S. $\langle mask \rangle$ the the $\langle mask \rangle$ $\langle mask \rangle$ $\langle mask \rangle$ tis00 culture R00 methods , she P $\langle mask \rangle$ the s U.00 $\langle mask \rangle$

We compare two methods

- BACKTRANS (NOISY) [Xie+2018]
 - Data is generated by back-translation
- DIRECTNOISE [Zhao+2019]
 - Data is generated by adding synthetic noise
- Please read our paper for details

Q3: How to use pseudo data during training?

"Genuine" Data



Q3: How to use pseudo data during training?



Experimental Configuration and Datasets

- We adapt "standard" configurations
 - Model: Transformer (Big) [Vaswani+2017]
 - Optimizer: Adam (for pretrain) and Adafactor (for finetuning)
- Dataset
 - BEA-2019 dataset (train/valid/test) [Bryant+2019]
 - CoNLL2014 (test) [Ng+2014]

Experiment 1: Choice for Seed Corpus

• Settings: JOINT

Method	Seed Corpus \mathcal{T}	Prec.	Rec.	$F_{0.5}$
Baseline	N/A	46.6	23.1	38.8

A1: Use Gigaword

- Seed corpus has minor influence on F_{0.5} score
- Gigaword is an ideal option
 - **clean text** is more important than **domain**?

Experiment 2: Utilization of Pseudo Data

• Settings: Wikipedia as seed corpus



If amount of pseudo data ≒ genuine data
 → PRETRAIN and JOINT are competitive

Experiment 2: Utilization of Pseudo Data

• Settings: Wikipedia as seed corpus



- Increasing amount of pseudo data improves the performance in PRETRAIN
- performance does not improve in JOINT
 - Pseudo data becomes dominant in JOINT

Experiment 3: More Pseudo Data

BACKTRANS (NOISY) significantly outperforms
 DIRECTNOISE







Comparison to Current Top Models

		CoNLL-2014 $(M^2 \text{ scorer})$			BEA-test (ERRANT)		
Model	Ensemble	Prec.	Rec.	$F_{0.5}$	Prec.	Rec.	$F_{0.5}$
Chollampatt and Ng (2018)	\checkmark	65.5	33.1	54.8	-	-	-
Lichtarge et al. (2019)	Previously Published Results 0.4						
Zhao et al. (2019) Grundkiewicz et al. (2019)	v √	/ 1.0	-		- 72.3	- 60.1	- 69.5
Grundatio (2017)	V			0112	12.5	00.1	07.0
LARGEPRETRAIN		67.9	44.1	61.3	65.5	59.4	64.2
LARGEPRETRAIN+SSE+R2L	LARGEPRET	RAIN	Results	5.0	72.1	61.8	69.8
LARGEPRETRAIN+SSE+R2L+S	SED 🗸	73.3	44.2	64.7	74.7	56.7	70.2

(1) Strong Single Model Results

		CoNLL-2014 $(M^2 \text{ scorer})$		BEA-test (ERRANT)			
Model	Ensemble	Prec.	Rec.	$F_{0.5}$	Prec.	Rec.	$F_{0.5}$
Chollampatt and Ng (2018)	\checkmark	65.5	33.1	54.8	-	-	-
Junczys-Dowmunt et al. (2018)	\checkmark	61.9	40.2	55.8	-	-	-
Lichtarge et al. (2019)	\checkmark	66.7	43.9	60.4	-	-	-
Zhao et al. (2019)	\checkmark	71.6	38.7	61.2	-	-	-
Grundkiewicz et al. (2019)	\checkmark	-	_(64.2	72.3	60.1	69.5
LARGEPRETRAIN		67.9	44.1	61.3	65.5	59.4	64.2
LARGEPRETRAIN+SSE+R2L	\checkmark	72.4	46.1	65.0	72.1	61.8	69.8
LARGEPRETRAIN+SSE+R2L+SED	\checkmark	73.3	44.2	64.7	74.7	56.7	70.2

Our best **single model** outperforms the **other ensemble models** except for [Grundkiewicz+2019]

(2) Additional Techniques Improve Result

			M^2 score	14 r)	BEA-test (ERRANT)		
Model	Ensemble	Prec.	Rec.	$F_{0.5}$	Prec.	Rec.	$F_{0.5}$
Chollampatt and Ng (2018)	\checkmark	65.5	33.1	54.8	-	-	-
Junczys-Dowmunt et al. (2018)	\checkmark	61.9	40.2	55.8	-	-	-
Lichtarge et al. (2019)	\checkmark	66.7	43.9	60.4	-	-	-
Zhao et al. (2019)	\checkmark	71.6	38.7	61.2	-	-	-
SSE: Synthetic Spelling Error	\checkmark	-	-	64.2	72.3	60.1	69.5
LARGEPF SED: Sentence-leve	LARGEPF SED: Sentence-level Error Det			61.3	65.5	59.4	64.2
LARGEPRETRAIN+SSE+R2L	\checkmark	72.4	46.1	65.0	72.1	61.8	69.8
LARGEPRETRAIN+SSE+R2L+SED	\checkmark	73.3	44.2	64.7	74.7	56.7	70.2
R2L: Right-to-Left Reranking Enser		e of 4 N	Nodels				
Our best model achieves t	he best pe	erform	ance c	on Co l	NLL20	14	

(F_{0.5}=65.0) and BEA-test (F_{0.5}=70.2)

Conclusions

• Investigated 3 questions regarding incorporating pseudo data into GEC model.

Q1: Choice for seed corpus? Q2: Methods for generating pseudo data? Q3: How to use pseudo data during training?

- Discovered settings suitable (LargePretrain)
 - justified by SOTA performance on benchmark datasets
- Code and pretrained model are available



November 18, 2019