

Competition or Cooperation? Exploring Unlabeled Data via Challenging Minimax Game for Semi-supervised Relation Extraction

Yu Hong¹, Jiahang Li¹, Jianchuan Feng¹, Chenghua Huang¹, Zhixu Li^{1*}, Jianfeng Qu², Yanghua Xiao^{1,3*}, Wei Wang¹

¹Shanghai Key Laboratory of Data Science, School of Computer Science, Fudan University

²School of Computer Science and Technology, Soochow University

³Fudan-Aishu Cognitive Intelligence Joint Research Center, Shanghai, China

{yhong17, jiahangli18, jcfeng20}@fudan.edu.cn, huangch22@m.fudan.edu.cn, zhixuli@fudan.edu.cn, jfqu@suda.edu.cn, {shawyh, weiwang1}@fudan.edu.cn

Abstract

Semi-Supervised Relation Extraction aims at learning well-performed RE models with limited labeled and large-scale unlabeled data. Existing methods mainly suffer from semantic drift and insufficient supervision, which severely limit the performance. To address these problems, recent work tends to design dual modules to work cooperatively for mutual enhancement. However, the consensus of two modules greatly restricts the model from exploring diverse relation expressions in unlabeled set, which hinders the performance as well as model generalization. To tackle this problem, in this paper, we propose a novel *competition-based* method *AdvSRE*. We set up a challenging minimax game on unlabeled data between two modules, Generator and Discriminator, and assign them with conflicting objectives. During the competition game, one module may find any possible chance to beat the other, which develops two modules' abilities until relation expressions cannot be further explored. To exploit label information, Discriminator is further asked to predict specific relation for each sentence. Experiment results on two benchmarks show new state-of-the-art performance over baselines, demonstrating the effectiveness of proposed AdvSRE.

Introduction

Relation Extraction (RE) plays an important role in natural language processing. It aims at extracting well-formed knowledge from large amounts of unstructured texts and has been widely used in many downstream tasks (Lin et al. 2019a; Wang et al. 2019; Shen et al. 2020). Given a sentence with two specified entities, the goal of RE is to identify the relation between two entities. For example, relation `the_writer_of` should be identified for the entity pair (`J.K.Rowling`, `Harry_Potter`) in the sentence "`J.K.Rowling` writes the much-loved series of `Harry_Potter` novels".

So far, plenty of neural RE methods are proposed, which provide end-to-end solutions and achieve promising performance in supervised RE (Zeng et al. 2014; Zhang et al. 2015; Vu et al. 2016). However, supervised neural models rely on large amounts of labeled data for effective training. Although distant supervision (Mintz et al. 2009) could allevi-

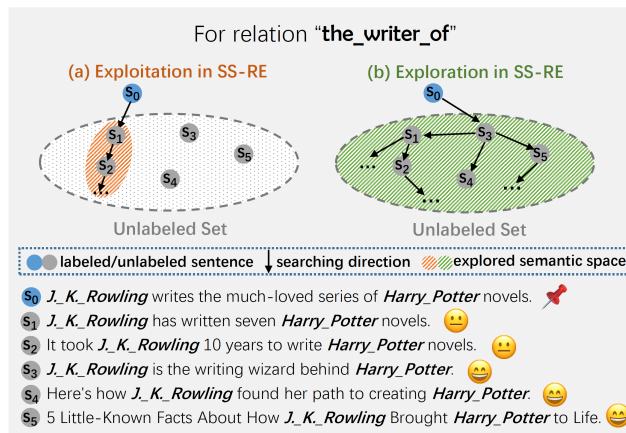


Figure 1: An example of exploitation and exploration of unlabeled sentences in SS-RE. Exploitation method tends to search sentences with similar expressions (e.g. s_1 and s_2) thus explores limited semantic space, while exploration method tends to search sentences with diverse relation expressions (e.g. s_3 - s_5) which can cover the whole space.

ate the requirement on manually labeling, it inevitably introduces noise. Given the drawbacks above on supervised and distantly-supervised RE, Semi-Supervised Relation Extraction (SS-RE) arises to learn well-performed RE models with limited labeled data and large amounts of unlabeled data (Agichtein and Gravano 2000; Sun and Grishman 2012).

Recent work in SS-RE tends to introduce two cooperative modules to alleviate semantic drift and insufficient supervision problems existing in conventional self-training (Paass 1993; Rosenberg, Hebert, and Schneiderman 2005) and self-ensembling (Tarvainen and Valpola 2017; Miyato et al. 2018) methods. For instance, DualRE (Lin et al. 2019b) takes sentence retrieval as a dual task of relation extraction, which trains the retrieval module jointly with the prediction module. Another work MetaSRE (Hu et al. 2021a) focuses on learning a generation network to generate high-quality pseudo labels for a relation classification network, which helps in return to meta-optimize the generation network. By making two modules benefit and correct each other, these so-called *cooperation-based* methods mitigate

*Corresponding authors.

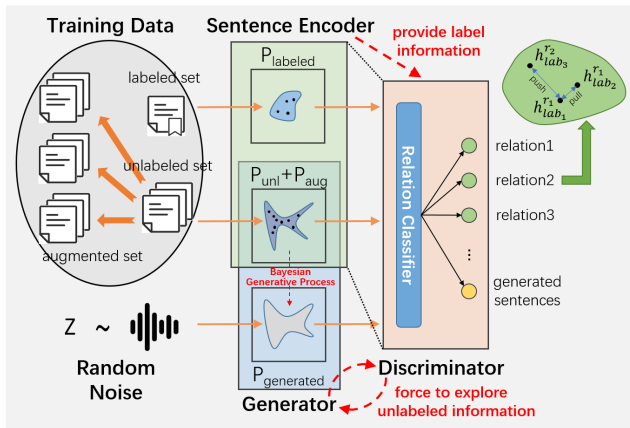


Figure 2: Overview of AdvSRE.

both insufficient supervision and semantic drift problems. However, the consensus of two modules excessively *exploits* acquired relation information from labeled set while restricts the model from *exploring* diverse relation expressions in unlabeled set. We illustrate this problem with an example in Figure 1 (a). Cooperation-based methods tend to exploit unlabeled data using the knowledge learned before, leaving other sentences with rich and diverse relation expressions unexplored. A recent method GradLRE (Hu et al. 2021b) tries to explore unlabeled set by rewarding positive annotations through reinforcement learning. However, to avoid the effect from wrong pseudo labels, it makes the gradient directions of unlabeled sentences imitate that learned from labeled data, which barely explores potential descent directions guided by unseen relation expressions. As a result, these methods present low recall and unsatisfactory performance, as reported in (Lin et al. 2019b; Hu et al. 2021a,b).

To tackle this problem, in this paper, we propose a novel *competition-based* method *AdvSRE*, which aims at fully exploring diverse relation expressions in unlabeled set to improve performance and generalization, as exemplified in Figure 1 (b). Specifically, we set up a challenging minimax game to make two modules *compete* with each other. By first augmenting sentences in unlabeled set in a label-preserving and diversity-enhanced manner, we make relation expressions in unlabeled set semantically rich. Then we set a module called *Generator* with powerful generating ability which is derived from data distribution in both unlabeled and augmented sets to produce high-quality fraudulent samples to fool the other module, *Discriminator*. Discriminator, on the other hand, is a powerful contextualized sentence encoder followed by a relation classifier which aims to correctly distinguish generated samples from real unlabeled ones. By making two powerful adversaries compete in this challenging game, Generator learns to capture rich data distribution in both unlabeled and augmented sets and generates highly-diverse samples to cheat, which forces Discriminator to fully explore the whole semantic space to win the game. When the competition game ends, both two modules are equipped with strong abilities to learn from diverse relation expressions in

unlabeled set, which is exactly what we need for exploration in SS-RE.

In addition to exploration on unlabeled set, we also exploit label information and propagate it to unlabeled set during the game. Specifically, when tackling with labeled sentences, we upgrade real-fake detection of Discriminator to *fine-grained relation extraction* and further pull sentences labeled with same relation together while push others apart. By effectively exploiting label information and fully exploring unlabeled information, we finally get a well-generalized model based on limited label data and large amounts of unlabeled data. The overview of proposed AdvSRE is shown in Figure 2. To summarize, our contributions are:

- We introduce *competition*, instead of cooperation mechanism in SS-RE to fully explore unlabeled sentences to improve model performance and generalization.
- We propose a framework *AdvSRE*, in which a *Generator* modeling distribution of unlabeled data and generating fake samples, a *Discriminator* capturing diverse relation expressions to not only distinguish real sentences from the fake but also act as the relation extractor.
- We develop *Bayesian Generative Process* for Generator to learn from wide distribution in unlabeled and augmented sets to generate highly-diverse fraudulent samples which help to inversely improve Discriminator.
- We upgrade real-fake detection of Discriminator to *fine-grained RE* and further learn *relational contrastive features* to effectively exploit label information to better propagate it to unlabeled set.
- We conduct extensive experiments on two benchmarks and achieve new state-of-the-art performance for SS-RE.

Related Work

Semi-Supervised Relation Extraction (SS-RE). There are roughly two lines of SS-RE methods, namely self-training and self-ensembling ones. Self-training methods (Paass 1993; Rosenberg, Hebert, and Schneiderman 2005) iteratively label a portion of unlabeled data and re-train the model based on the updated training set. However, they inevitably introduce wrong pseudo labels thus accumulate learning errors, known as semantic drift (Curran, Murphy, and Scholz 2007). Self-ensembling methods (Tavainen and Valpola 2017; Miyato et al. 2018) come from the idea that better performance could be obtained by ensembling models under perturbations of data. However, they heavily rely on limited supervision from labeled data and suffer from insufficient supervision. To alleviate these problems, recent work on SS-RE tends to design two modules to work cooperatively for mutual enhancement. RE-Ensemble (Lin et al. 2019b) proposes two independent modules to generate pseudo labels and select the intersection between their predictions as high-quality annotations. DualRE (Lin et al. 2019b) further replaces one module as a retrieval module which helps to improve the other prediction module. MetaSRE (Hu et al. 2021a) adopts meta learning for the generation module to ensure the quality of pseudo labels. However, the mode of cooperation restricts two modules

from exploring diverse relation expressions, which limits the model from generalising beyond knowledge learned before. GradLRE (Hu et al. 2021b) tries to explore unlabeled set via reinforcement learning. It quantifies annotating behavior and rewards positive feedbacks during trial and error. However, GradLRE relies on gradient imitation of labeled data to avoid wrong annotations on unlabeled sentences, which may easily fall into local minima when searching gradient descent directions. Proposed AdvSRE, on the other hand, introduces competition between two modules to enforce full exploration throughout the whole unlabeled set. It assigns no labels for unlabeled data and also avoids semantic drift and insufficient supervision due to adversarial learning.

Generative Adversarial Learning (GAL). GAL (Goodfellow et al. 2014) provides a framework to estimate data distribution with two adversarial components. The quality of the generated samples can be as indistinguishable as real ones. GAL has been widely used in CV as well as NLP tasks, such as face synthesis (Liu, Li, and Sun 2019; Fu et al. 2019) and dialogue generation (Su et al. 2018; Zhu et al. 2021). For relation extraction, (Qin, Xu, and Wang 2018) uses GAL to filter noisy samples in distantly-supervised training set. (Li et al. 2019) constructs clean sets based on knowledge base then makes use of NA sentences with GAL. (Luo, Pan, and Peng 2020) learns the distribution of true positive instances then generates valid sentences for model training. (Hao, Yu, and Hu 2021) aligns filtered false negative sentences with positive ones then redistributes them to real relations. These works are all based on distantly-supervised RE. To the best of our knowledge, we are the first to extend GAL to SS-RE.

Contrastive Learning (CL). CL is widely used in various tasks to help to improve latent representations. (Yan et al. 2021) uses CL to avoid collapsed sentence representations learned from BERT. (Wang et al. 2021) studies the impact on representations under different perturbations. (Gunel et al. 2021) improves the robustness on few-shot tasks. We leverage the method in (Gunel et al. 2021) to improve the performance under semi-supervised data settings.

Proposed Method AdvSRE

In AdvSRE, there are two adversarial modules: (1) a Generator G , which transforms noise vector into fraudulent sentence representation h_{gen} ; (2) a Discriminator D , which includes a Sentence Encoder to map real sentence in labeled, unlabeled and augmented sets into latent representation h_{lab} , h_{unl} and h_{aug} , and a Relation Classifier to detect relation \hat{r} for h_{gen} , h_{lab} , h_{unl} and h_{aug} , respectively. During the competition game, G learns to generate high-quality fake samples from unlabeled set to extend D 's ability on exploring diverse relation expressions in semantic space. Meanwhile, D also learns to correctly classify labeled sentences into their corresponding relations, leading to fine-grained relation extraction throughout the whole training process. We illustrate these two modules as follows.

Generator

Generator in AdvSRE aims to imitate unlabeled sentences to generate fraudulent samples to improve the Discriminator.

Original s : As we saw earlier, *helicobacter* is responsible for causing stomach *ulcer*.

Back-translated s' : As we saw earlier, *helicobacter* is the cause of *ulcer* in the stomach.

TF-IDF replaced s^* : As we saw earlier, *helicobacter* is responsible for beta stomach *ulcer*.

Table 1: Augmented sentences using back-translation and TF-IDF replacement.

Since the diversity in fake samples is the key to promote D 's ability on identifying different kinds of relation expressions, we take two measures to ensure this: the first one is to augment unlabeled sentences in a label-preserving way, which means we enrich the relation expressions between two entities but do not alter their relations. The second one is to assign Generator with powerful ability to learn from the whole distribution in both unlabeled and augmented sets, make it generate diverse and high-quality fraudulent samples for deceit. These two measures make the competition game between G and D challenging, since D is forced to explore the whole feature space to identify fake samples to beat G . In SS-RE, this is exactly what we need to promote D as a powerful and generalized relation extractor.

Augmentation for unlabeled sentences. We adopt back-translation (Edunov et al. 2018) and TF-IDF replacement (Xie et al. 2020; Chen et al. 2021) as our augmentation methods, since the former introduces diversity by reformulating the whole sentence while the latter locally changes some words. Moreover, they are both unsupervised methods which do not change the relation between two entities.

- **Back-translation.** We translate unlabeled sentence s into Chinese/French, then translate it back again to obtain augmented sentence s' . To maintain entity mentions in s' , we replace them with special tokens to avoid being substituted during augmentation.
- **TF-IDF replacement.** We consider unlabeled sentence s as a document and the whole unlabeled set U as the corpus, calculate TF-IDF score for each word then replace unimportant words to form s^* . We keep entity mentions in s^* unchanged as in s' .

An example of sentence augmentation using two methods above can be found in Table 1. The augmented samples $\{s', s^*\}$ for each s from U form the augmented set A .

Learning from unlabeled data's distribution. To fully capture data distribution in both unlabeled and augmented sets, we put a distribution over parameters of a two-layer perceptron as our Generator G . Compared with G with fixed parameters, this method results in a bunch of generators whose parameters are drawn from the broad distribution of rich expressions in unlabeled and augmented sentences. With variously different generators, we can produce highly-diverse fraudulent samples from the whole semantic space and also ensure they are as close as unlabeled or augmented sentences. Formally, we put a prior $P(\theta_g)$ over G 's param-

Algorithm 1: Inference of Bayesian Generative Process

Input: Labeled set T , unlabeled set U , augmented set A , generated set F

Parameter: Learning rate η , friction term γ , number of MC iterations K , number of SGHMC updates L

Output: G 's posterior represented by sample set Θ_g , D 's parameter θ_d

```
1:  $\Theta_g \leftarrow \emptyset$ 
2: for  $k := 1, \dots, K$  do
3:   for  $l := 1, \dots, L$  do
4:      $q \sim N(0, 2\gamma\eta I)$ 
5:      $u \leftarrow (1 - \gamma)u + q + \eta \nabla_{\theta_g} \log[P(\theta_g|F, \theta_d)]$ 
6:      $\theta_g \leftarrow \theta_g + u$ 
7:   end for
8:    $\Theta_g \leftarrow \Theta_g \cup \theta_g$ 
9: end for
10:  $\theta_d \leftarrow \theta_d + \eta \nabla_{\theta_d} \log[P(\theta_d|T, U, A, F, \Theta_g)]$ 
11: return  $\Theta_g, \theta_d$ 
```

ter θ_g , and sample from its posterior $P(\theta_g|U, A)$ as:

$$\theta_g \sim P(\theta_g|U, A) \quad (1)$$

Then generated representation h_{gen} is derived from G with sampled parameter θ_g and the noise vector Z :

$$h_{gen} = G(Z; \theta_g), Z \sim N(0, I) \quad (2)$$

where $Z \in \mathbb{R}^{100}$ and $h_{gen} \in \mathbb{R}^{2d}$. h_{gen} has the same dimension with real sentence's representation which will be illustrated in next section.

To get h_{gen} in practice, we first (1) draw a value of θ_g from $P(\theta_g|U, A)$, then (2) draw m different noise vector Z from $N(0, I)$, finally (3) condition G on parameter θ_g to transform each Z into h_{gen} . We do this process K times to get mK generated samples for each unlabeled and augmented sentence as our high-quality fraudulent set F . We denote this method as Bayesian Generative Process (**BGP**).

In practice, G is activated by Leaky-ReLU and dropped out with a certain rate. To sample G 's parameters, we adopt Stochastic Gradient Hamiltonian Monte Carlo (SGHMC) (Chen, Fox, and Guestrin 2014) and meanwhile inference D 's parameters during the competition game. One iteration of BGP's inference is presented in Algorithm 1. BGP can be considered as introducing useful inductive bias (Saatchi and Wilson 2017). It makes G produce more different fake representations to push D to be a better relation extractor.

Discriminator

In order to compete against Generator in the challenging game, we also equip Discriminator with powerful ability on extracting contextualized features from sentences and identifying distinguishable features from different relations. It consists of two components, i.e. Sentence Encoder and Relation Classifier, which are explained as follows.

Sentence Encoder. To extract rich context information of entity pairs from labeled, unlabeled and augmented sentences, we adopt BERT (Devlin et al. 2018) as our

sentence encoder. To further capture entity-level features for effective relation extraction, we adopt the tagging schema in (Soares et al. 2019) to insert four symbols $\langle e1 \rangle$, $\langle /e1 \rangle$, $\langle e2 \rangle$, and $\langle /e2 \rangle$ around two entity mentions. [CLS] and [SEP] are also added at the beginning and the end of the sentence. For example, a sentence (i.e. s_{lab} , s , s' or s^*) with words $\{w_1, \dots, w_n\}$ becomes $\{[CLS], w_1, \dots, \langle e1 \rangle, w_{e1}, \langle /e1 \rangle, \dots, \dots, \langle e2 \rangle, w_{e2}, \langle /e2 \rangle, \dots, w_n, [SEP]\}$.

Instead of using the representation of [CLS] as sentence-level features, we concatenate BERT's outputs of $\langle e1 \rangle$ and $\langle e2 \rangle$ as entity-level features. Let $h_{\langle e1 \rangle}$ and $h_{\langle e2 \rangle}$ denote the hidden representations of $\langle e1 \rangle$ and $\langle e2 \rangle$, entity-level representation of a real sentence h_{real} is given by:

$$h_{real} = [h_{\langle e1 \rangle}; h_{\langle e2 \rangle}] \quad (3)$$

where $h_{real} \in \mathbb{R}^{2d}$, d is the dimension of BERT's output. $h_{real} = \{h_{lab}, h_{unl}, h_{aug}\}$.

Relation Classifier. To detect corresponding relation for each sample's representation, we implement Relation Classifier as a two-layer perceptron activated by Leaky-ReLU and dropped out with a certain rate to transform sentence's representation into logit vector. Then we conduct Softmax on the logits and select the relation with maximum probability as the predictive relation \hat{r} :

$$\hat{r} = \arg \max_r [\text{Softmax}(\text{MLP}(h))] \quad (4)$$

where $h = \{h_{real}, h_{gen}\}$.

With powerful Sentence Encoder and Relation Classifier, D has the ability to compete against G to identify fake samples from real sentences. Moreover, we extend D 's real-fake detection from binary classification to multi-classification (which is embodied in loss functions and will be explained in next section), resulting in fine-grained RE during the whole training process. If labeled data traverse all subdomains of the feature space, which is a reasonable assumption under semi-supervised conditions, relation information can then be propagated from labeled set to the whole feature space.

Loss Functions

Since our aim is to generate fraudulent samples from Generator to force Discriminator to explore unlabeled (and also augmented) sets, the optimization objectives on generated set are adversarial for two modules. For Generator G , we aim to fool D to misclassify these samples into real relations, so we minimize the loss of generated set F on top R real classes:

$$L_{G_{gen}} = -\mathbb{E}_{h_{gen} \in F} [\log \sum_{r=1}^R P(\hat{r} = r | h_{gen})] \quad (5)$$

For Discriminator, on the contrary, we expect generated samples to be classified into $(R + 1)$ -th fake class. This indicates D has the ability to fully explore the feature space to distinguish real sentences from the fake. Loss function of D on fraudulent set F is defined as:

$$L_{D_{gen}} = -\mathbb{E}_{h_{gen} \in F} [\log P(\hat{r} = R + 1 | h_{gen})] \quad (6)$$

¹For simplicity, we omit the loss term on G 's prior $N(0, I)$.

Methods	% Labeled Data			5%			10%			30%		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Self-Training (Rosenberg, Hebert, and Schneiderman 2005)	74.30	71.78	73.02	76.64	74.10	75.35	81.92	83.09	82.50			
Mean-Teacher (Tarvainen and Valpola 2017)	73.23	72.65	72.94	75.43	73.18	74.79	79.69	83.23	81.42			
RE-Ensemble (Lin et al. 2019b)	73.88	70.88	72.35	76.83	74.62	75.71	81.26	81.42	81.34			
DualRE (Lin et al. 2019b)	74.12	78.21	76.11	76.71	79.81	78.23	82.10	85.05	83.55			
MRefG (Li et al. 2021)	73.04	78.29	75.48	76.32	79.76	77.96	81.75	84.91	83.24			
MetaSRE (Hu et al. 2021a)	75.59	81.40	78.33	78.05	82.29	80.09	82.01	87.95	84.81			
GradLRE (Hu et al. 2021b)	76.62	81.62	79.65	78.99	84.58	81.69	83.84	87.27	85.52			
AdvSRE (ours)	78.23	85.94	81.90	79.16	86.61	82.72	84.47	88.47	86.42			
Fully-Supervised RE (Soares et al. 2019)	84.15	85.14	84.64	84.37	86.46	85.40	86.51	88.13	87.08			

Table 2: Performance on SemEval with different proportions of labeled data and 50% unlabeled data.

Methods	% Labeled Data			3%			10%			15%		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Self-Training (Rosenberg, Hebert, and Schneiderman 2005)	49.12	38.47	43.15	56.91	52.66	54.70	60.02	54.12	56.92			
Mean-Teacher (Tarvainen and Valpola 2017)	53.12	40.75	46.12	58.12	50.58	54.09	58.00	52.60	55.17			
RE-Ensemble (Lin et al. 2019b)	51.39	36.64	42.78	57.34	52.53	54.83	61.19	51.08	55.68			
DualRE (Lin et al. 2019b)	59.23	36.01	44.79	60.92	52.82	56.58	61.48	56.09	58.66			
MRefG (Li et al. 2021)	56.31	36.25	43.81	59.25	51.93	55.42	61.02	55.61	58.21			
MetaSRE (Hu et al. 2021a)	58.96	37.66	46.16	60.49	53.69	56.95	65.03	54.02	58.94			
GradLRE (Hu et al. 2021b)	54.67	41.79	47.37	62.26	54.64	58.20	65.32	55.36	59.93			
AdvSRE (ours)	55.97	42.60	48.38	60.73	61.48	61.10	61.07	64.05	62.53			
Fully-Supervised RE (Soares et al. 2019)	66.27	60.66	63.34	67.54	60.61	63.89	68.32	61.95	64.98			

Table 3: Performance on TACRED with different proportions of labeled data and 50% unlabeled data.

For unlabeled set U , D tries to classify the sentences into top R real relations to distinguish them from fake samples. To do so, we maximize the sum of D 's probabilities on top R classes. Loss function of D on unlabeled set U is:

$$L_{D_{unl}} = -\mathbb{E}_{h_{unl} \in U} [\log \sum_{r=1}^R P(\hat{r} = r | h_{unl})] \quad (7)$$

Augmented set can be considered as the diverse twin of unlabeled set. For augmented set A , we optimize the loss on first R relations like what we do on unlabeled set:

$$L_{D_{aug}} = -\mathbb{E}_{h_{aug} \in A} [\log \sum_{r=1}^R P(\hat{r} = r | h_{aug})] \quad (8)$$

For labeled set T , D should correctly identify specific relation for each labeled sentence to propagate label information to the whole feature space. Loss function of D on labeled set T is defined as the cross-entropy between labeled sentence's representation and its corresponding relation r :

$$L_{D_{lab}} = -\mathbb{E}_{h_{lab} \in T} [\log P(\hat{r} = r | h_{lab})] \quad (9)$$

To better exploit label information, we guide Discriminator to further learn Relational Contrastive Features (RCF) from labeled sentences. Following (Gunel et al. 2021), we pull sentences expressing same relation together while push other sentences apart. The loss of RCF is defined as:

$$L_{D_{RCF}} = -\mathbb{E}_{h_{lab} \in T} [\mathbb{E}_{h'_{lab} \in T, h'_{lab} \neq h_{lab}, r' = r} [\log \frac{\exp(\cos(h_{lab}, h'_{lab})/\tau)}{\sum_{h''_{lab} \in T, h''_{lab} \neq h_{lab}} \exp(\cos(h_{lab}, h''_{lab})/\tau)}]] \quad (10)$$

where τ is the scaled factor, $\cos(\cdot, \cdot)$ represents the cosine similarity between two sentence representations.

Total loss of D is defined as the sum of loss on generated, unlabeled, augmented, labeled sets and also the loss on RCF:

$$L_D = L_{D_{gen}} + L_{D_{unl}} + L_{D_{aug}} + (1 - \beta)L_{D_{lab}} + \beta L_{D_{RCF}} \quad (11)$$

where β is the coefficient parameter which makes the balance between $L_{D_{lab}}$ and $L_{D_{RCF}}$.

Experiments and Analysis

In this section, we first introduce datasets, baselines and experimental settings for SS-RE, then we present performance of baselines and proposed AdvSRE with detailed analysis.

Datasets

We conduct our experiments on two standard benchmarks:

- **SemEval**² (Hendrickx et al. 2010) is a popular RE dataset. It contains 7,199 sentences in training set, 800 sentences in validation set and 1,864 sentences in test set. It has 19 relations in total, including `no_relation` which indicates there is no relation between two entities. The proportion of `no_relation` sentences is 17.4%.
- **TACRED**³ (Zhang et al. 2017) is a larger benchmark with 75,049 sentences in training set, 25,763 sen-

²<http://semeval2.fbk.eu/semeval2.php>

³<https://catalog.ldc.upenn.edu/LDC2018T24>

tences in validation set and 18,659 sentences in test set. It is more complicated with 42 relations (including `no_relation`) and more skewed with 78.68% `no_relation` sentences.

Baselines

We adopt pre-trained language model BERT (Devlin et al. 2018) as the encoder and compare with 8 strong baselines:

- **Self-Training** (Rosenberg, Hebert, and Schneiderman 2005) iteratively trains the model on labeled set and generates pseudo labels on unlabeled set. It stops when unlabeled data is exhausted.
- **Mean-Teacher** (Tarvainen and Valpola 2017) is a self-ensembling method which gathers outputs from models under parameter perturbations.
- **RE-Ensemble** (Lin et al. 2019b) uses two independent prediction modules to infer relations for unlabeled sentences and select pseudo labels according to the agreement on their prediction results.
- **DualRE**⁴ (Lin et al. 2019b) considers sentence retrieval as the dual problem of relation extraction. It contains a retrieval module and an extraction module and train them jointly to promote newly-generated pseudo labels.
- **MRefG** (Li et al. 2021) correlates unlabeled sentences with labeled sentences by constructing entity, verb and semantics reference graphs.
- **MetaSRE**⁵ (Hu et al. 2021a) adopts meta-learning at the beginning of training to generate high-quality pseudo labels for unlabeled sentences.
- **GradLRE**⁶ (Hu et al. 2021b) is the current state-of-the-art SS-RE method. It adopts reinforcement learning to reward the annotations that imitate the behavior of labeled data on model gradient descent directions.
- **Fully-Supervised RE** (Soares et al. 2019) is adopted as the performance ceiling of SS-RE. It is fully-supervised by all attainable training data (i.e. both labeled and unlabeled sentences with golden labels).

Experimental Settings

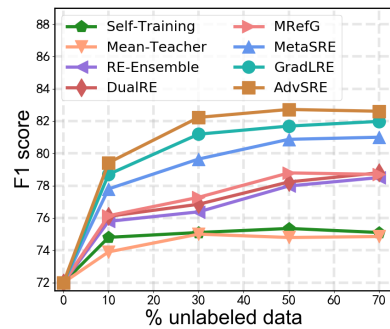
Following (Rosenberg, Hebert, and Schneiderman 2005; Tarvainen and Valpola 2017; Lin et al. 2019b; Li et al. 2021; Hu et al. 2021a,b), we adopt F1 score as the evaluation metric and precision and recall as auxiliary metrics. For data settings, we follow (Lin et al. 2019b; Hu et al. 2021a,b) to divide the training set into labeled and unlabeled sets. For SemEval, we sample 5%, 10% and 30% of original training set as labeled sets and 50% as unlabeled set. For TACRED, we sample 3%, 10% and 15% as labeled sets and 50% as unlabeled set. We adopt stratified sampling in (Lin et al. 2019b; Hu et al. 2021a,b) to ensure relation proportion does not change in both labeled and unlabeled sets.

For parameter settings, we set sentence’s maximum length as 128 and batch size as 16. We adopt AdamW as

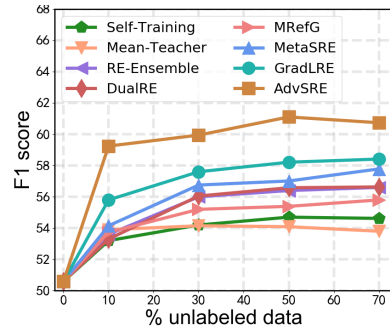
⁴<https://github.com/INK-USC/DualRE>

⁵<https://github.com/THU-BPM/MetaSRE>

⁶<https://github.com/THU-BPM/GradLRE>



(a) F1 scores on SemEval.



(b) F1 scores on TACRED.

Figure 3: Performance on SemEval and TACRED with different proportions of unlabeled data and 10% labeled data.

the optimizer and set learning rate as $5e - 5$. We warm-up the learning rate in the first 50 steps then linearly decrease it to 0. For BGP in AdvSRE, we set friction term of SGHMC (i.e. γ) as 0.001, number of MC iterations (i.e. K) as 10 and number of SGHMC updates (i.e. L) as 1. For RCF, we set $\tau = 0.3$, $\beta = 0.1$. Dropout rate is set as 0.1 for both Discriminator and Generator. Following (Hu et al. 2021a,b), we set the training epoch as 10 and run 5 times training and testing to report the average performance. For baselines, we adopt the parameter settings in the original papers.

Results and Analysis

In this section, we analyze both quantitative and qualitative effectiveness of proposed method in terms of performance, feature exploration, model generalization and different competition strategies.

Performance on different proportions of labeled data.

Table 2 and 3 show the performance on SemEval and TACRED with different proportions of labeled and a fixed proportion of unlabeled data. As labeled data increases, we observe that AdvSRE consistently outperforms all baselines on all data settings. Specifically, AdvSRE gets average improvement of 1.71% F1 score on SemEval and 3.82% F1 score on TACRED compared with GradLRE, achieving new state-of-the-art performance in SS-RE. When compared with Fully-Supervised RE, performance gap is further narrowed, especially when more labeled data is provided. We also observe that AdvSRE can achieve much higher recall

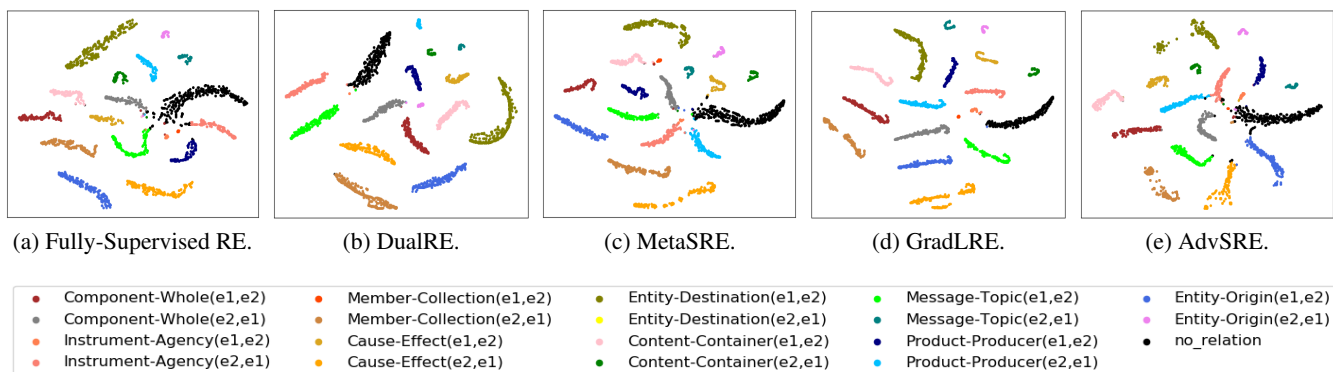
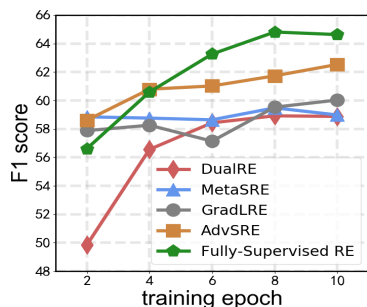
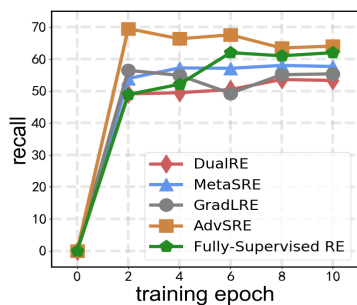


Figure 4: Exploration on 50% unlabeled data of SemEval with different methods using 10% labeled data.



(a) F1 score on different epochs.



(b) Recall on different epochs.

Figure 5: F1 score and recall of different methods on different training epochs on TACRED with 15% labeled and 50% unlabeled data.

than Fully-Supervised RE. This can be attributed to the challenging competition game set up on unlabeled sentences, which forces two modules to fully explore diverse relation expressions, thus extensively improves generalization.

Performance on different proportions of unlabeled data.

Figure 3 shows the performance on SemEval and TACRED with different proportions of unlabeled and a fixed proportion of labeled data. With the increase of unlabeled data, we can see that all methods get improvements on F1 score. Among them, AdvSRE achieves best performance and out-

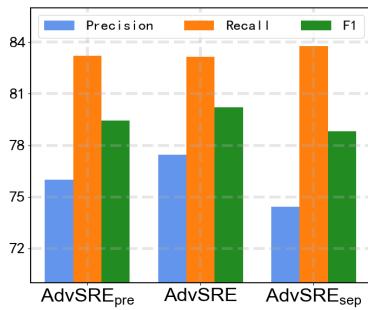
performs all baselines with a large margin (especially on TACRED dataset), demonstrating the effectiveness of competition mechanism in SS-RE.

Performance on feature exploration. Figure 4 shows feature exploration in semantic space for 50% unlabeled sentences in SemEval. To do so, we obtain unlabeled sentences' representations when model training is finished, then reduce the dimension with t-SNE (Hinton 2008). In Figure 4, we observe that all 5 methods get clear boundary for each relation. For DualRE, MetaSRE and GradLRE, sentences' representations distribute densely, since they tend to exploit feature space according to acquired knowledge of relation expressions. While for AdvSRE, sentences in the same relation seem to be more scattered. This is because we try to explore diverse relation expressions during the challenging competition game, resulting in broader exploration in the whole semantic space.

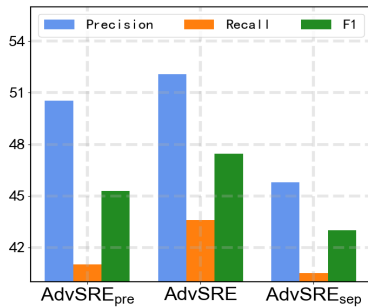
Performance on model generalization. To present model generalization, we train AdvSRE and baselines with 15% labeled and 50% unlabeled data of TACRED then report performance on each epoch, as shown in Figure 5. In Figure 5(a), as training epoch grows, performance of all methods is improved. Among them AdvSRE consistently outperforms all baselines, ending up at 62.53 at the last epoch, much higher than the current SOTA GradLRE. In Figure 5(b), AdvSRE shows best generalization at the beginning of training, reaching the recall of 70.52 and is much higher than Fully-Supervised RE. As training process goes, AdvSRE keeps steady on the highest recall, indicating best generalization in different training process of SS-RE.

Performance under different competition strategies.

To prove the effectiveness of proposed competition strategy, we design two other methods for the challenging game. The first one is to pre-train D on labeled data then simultaneously train it with G on both labeled and unlabeled sets, which is denoted as AdvSRE_{pre}. The second one is to simplify the training process by separately training D on labeled data and making it compete with G on unlabeled data, which is denoted as AdvSRE_{sep}. Results of AdvSRE_{pre}, AdvSRE and AdvSRE_{sep} with 5%, 3% labeled and 50% unlabeled



(a) Performance on SemEval.



(b) Performance on TACRED.

Figure 6: Performance of different competition strategies on SemEval and TACRED with 5%, 3% labeled and 50% unlabeled data.

data are shown in Figure 6. We find that AdvSRE_{pre} does not help too much for performance improvement. It even causes overfitting as precision and recall both decrease on two datasets. AdvSRE_{sep} underperforms and presents much lower performance on TACRED. Instead, AdvSRE conducts feature exploration and fine-grained RE at the same time, leading to the best performance among all strategies.

Conclusion

In this paper, we propose competition against cooperation in SS-RE by setting up a challenging minimax game between two modules to fully explore diverse relation expressions in unlabeled set. We also exploit label information by fine-grained RE with relational contrastive features. Experiment results show new state-of-the-art performance for SS-RE.

Acknowledgements

This work is supported by National Key Research and Development Project (No.2020AAA0109302), Shanghai Municipal Science and Technology Major Project (No.2021SHZDZX0103), Science and Technology Commission of Shanghai Municipality Grant (No.22511105902), National Natural Science Foundation of China (No.62102095; No.62072323; No.62102276), Shanghai Science and Technology Innovation Action Plan (No.22511104700), Natural Science Foundation of Jiangsu

Province (No.BK20210705), and Natural Science Foundation of Educational Commission of Jiangsu Province, China (No.21KJD520005). Yanghua Xiao is also a member of Research Group of Computational and AI Communication at Institute for Global Communications and Integrated Media, Fudan University.

References

- Agichtein, E.; and Gravano, L. 2000. Snowball: Extracting relations from large plain-text collections. In *Proceedings of the fifth ACM conference on Digital libraries*, 85–94.
- Chen, T.; Fox, E.; and Guestrin, C. 2014. Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, 1683–1691. PMLR.
- Chen, T.; Shi, H.; Tang, S.; Chen, Z.; Wu, F.; and Zhuang, Y. 2021. CIL: Contrastive Instance Learning Framework for Distantly Supervised Relation Extraction. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 6191–6200. Association for Computational Linguistics.
- Curran, J. R.; Murphy, T.; and Scholz, B. 2007. Minimising semantic drift with mutual exclusion bootstrapping. In *Proceedings of the 10th Conference of the Pacific Association for Computational Linguistics*, volume 6, 172–180. Citeseer.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Edunov, S.; Ott, M.; Auli, M.; and Grangier, D. 2018. Understanding back-translation at scale. *arXiv preprint arXiv:1808.09381*.
- Fu, C.; Wu, X.; Hu, Y.; Huang, H.; and He, R. 2019. Dual variational generation for low shot heterogeneous face recognition. *Advances in Neural Information Processing Systems*, 32.
- Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. *arXiv preprint arXiv:1406.2661*.
- Gunel, B.; Du, J.; Conneau, A.; and Stoyanov, V. 2021. Supervised contrastive learning for pre-trained language model fine-tuning. In *International Conference on Learning Representations*.
- Hao, K.; Yu, B.; and Hu, W. 2021. Knowing False Negatives: An Adversarial Training Method for Distantly Supervised Relation Extraction. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 9661–9672.
- Hendrickx, I.; Kim, S. N.; Kozareva, Z.; Nakov, P.; Séaghdha, D. O.; Padó, S.; Pennacchiotti, M.; Romano, L.; and Szpakowicz, S. 2010. Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. *arXiv preprint arXiv:1911.10422*.
- Hinton, G. 2008. Visualizing High-Dimensional Data Using t-SNE. *Journal of Machine Learning Research*, 9: 2579–2605.

- Hu, X.; Ma, F.; Liu, C.; Zhang, C.; Wen, L.; and Yu, P. S. 2021a. Semi-supervised Relation Extraction via Incremental Meta Self-Training. In *Findings of EMNLP*.
- Hu, X.; Zhang, C.; Yang, Y.; Li, X.; Lin, L.; Wen, L.; and Yu, P. S. 2021b. Gradient Imitation Reinforcement Learning for Low Resource Relation Extraction. In *EMNLP*.
- Li, P.; Zhang, X.; Jia, W.; and Zhao, H. 2019. GAN driven semi-distant supervision for relation extraction. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 3026–3035.
- Li, W.; Qian, T.; Chen, X.; Tang, K.; Zhan, S.; and Zhan, T. 2021. Exploit a Multi-head Reference Graph for Semi-supervised Relation Extraction. In *IJCNN*.
- Lin, B. Y.; Chen, X.; Chen, J.; and Ren, X. 2019a. Kagnet: Knowledge-aware graph networks for commonsense reasoning. *arXiv preprint arXiv:1909.02151*.
- Lin, H.; Yan, J.; Qu, M.; and Ren, X. 2019b. Learning dual retrieval module for semi-supervised relation extraction. In *The World Wide Web Conference*, 1073–1083.
- Liu, Y.; Li, Q.; and Sun, Z. 2019. Attribute-aware face aging with wavelet-based generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11877–11886.
- Luo, G.; Pan, J.; and Peng, M. 2020. RDSGAN: Rank-based distant supervision relation extraction with generative adversarial framework. *arXiv preprint arXiv:2009.14722*.
- Mintz, M.; Bills, S.; Snow, R.; and Jurafsky, D. 2009. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, 1003–1011.
- Miyato, T.; Maeda, S.-i.; Koyama, M.; and Ishii, S. 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8): 1979–1993.
- Paass, G. 1993. Assessing and improving neural network predictions by the bootstrap algorithm. In *Advances in Neural Information Processing Systems*, 196–203. Citeseer.
- Qin, P.; Xu, W.; and Wang, W. Y. 2018. DSGAN: Generative Adversarial Training for Distant Supervision Relation Extraction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 496–505.
- Rosenberg, C.; Hebert, M.; and Schneiderman, H. 2005. Semi-Supervised Self-Training of Object Detection Models. In *Proceedings of the Seventh IEEE Workshops on Application of Computer Vision (WACV/MOTION'05)-Volume 1-Volume 01*, 29–36.
- Saatchi, Y.; and Wilson, A. G. 2017. Bayesian gan. *arXiv preprint arXiv:1705.09558*.
- Shen, T.; Mao, Y.; He, P.; Long, G.; Trischler, A.; and Chen, W. 2020. Exploiting structured knowledge in text via graph-guided representation learning. *arXiv preprint arXiv:2004.14224*.
- Soares, L. B.; Fitzgerald, N.; Ling, J.; and Kwiatkowski, T. 2019. Matching the Blanks: Distributional Similarity for Relation Learning. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2895–2905.
- Su, H.; Shen, X.; Hu, P.; Li, W.; and Chen, Y. 2018. Dialogue generation with GAN. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Sun, A.; and Grishman, R. 2012. Active learning for relation type extension with local and global data views. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, 1105–1112.
- Tarvainen, A.; and Valpola, H. 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30.
- Vu, N. T.; Adel, H.; Gupta, P.; and Schütze, H. 2016. Combining recurrent and convolutional neural networks for relation classification. *arXiv preprint arXiv:1605.07333*.
- Wang, D.; Ding, N.; Li, P.; and Zheng, H.-T. 2021. Cline: Contrastive learning with semantic negative examples for natural language understanding. *arXiv preprint arXiv:2107.00440*.
- Wang, X.; He, X.; Cao, Y.; Liu, M.; and Chua, T.-S. 2019. Kgat: Knowledge graph attention network for recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 950–958.
- Xie, Q.; Dai, Z.; Hovy, E.; Luong, T.; and Le, Q. 2020. Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems*, 33: 6256–6268.
- Yan, Y.; Li, R.; Wang, S.; Zhang, F.; Wu, W.; and Xu, W. 2021. Consert: A contrastive framework for self-supervised sentence representation transfer. *arXiv preprint arXiv:2105.11741*.
- Zeng, D.; Liu, K.; Lai, S.; Zhou, G.; and Zhao, J. 2014. Relation classification via convolutional deep neural network. In *Proceedings of COLING 2014, the 25th international conference on computational linguistics: technical papers*, 2335–2344.
- Zhang, S.; Zheng, D.; Hu, X.; and Yang, M. 2015. Bidirectional long short-term memory networks for relation classification. In *Proceedings of the 29th Pacific Asia conference on language, information and computation*, 73–78.
- Zhang, Y.; Zhong, V.; Chen, D.; Angeli, G.; and Manning, C. D. 2017. Position-aware attention and supervised data improve slot filling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 35–45.
- Zhu, Q.; Chen, X.; Wu, P.; Liu, J.; and Zhao, D. 2021. Combining Curriculum Learning and Knowledge Distillation for Dialogue Generation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 1284–1295.