

Minimally-Supervised Relation Induction from Pre-trained Language Model

Anonymous ACL submission

Abstract

Relation Induction is a very practical task in Natural Language Processing (NLP) area. In practical application scenarios, people want to induce more entity pairs having the same relation from only a few seed entity pairs. Thus, instead of the laborious supervised setting, in this paper, we focus on the minimally-supervised setting where only a couple of seed entity pairs per relation are provided. Although the conventional relation induction methods have made some success, their performance depends heavily on the quality of word embeddings. The great success of Pre-trained Language Models, such as BERT, changes the NLP area a lot, and they are proven to be able to better capture relation knowledge. In this paper, we propose a novel method to induce relation with BERT under the minimally-supervised setting. Specifically, we firstly extract proper templates from the corpus by using the mask-prediction task in BERT to build pseudo-sentences as the context of entity pairs. Then we use BERT attention weights to better represent the pseudo-sentences. In addition, We also use the Integrated Gradient of entity pairs to iteratively select better templates further. Finally, with the high-quality pseudo-sentences, we can train a better classifier for relation induction. Experiments on Google Analogy Test Sets (GATS), Bigger Analogy Test Set (BATS) and DiffVec demonstrate that our proposed method achieves state-of-the-art performance.

1 Introduction

Relation induction is a task to judge whether two entities have a certain relation based on some given entity pairs of that relation, which was first proposed in (Vylomova et al., 2016). For instance, given $\{(Germany, Berlin), (France, Paris), (Italy, Rome)\}$, relation induction is to predict whether new entity pairs such as $(China, Beijing)$ have the same relation as the given entity pairs. In practical scenarios, only a

few seed entity pairs are available. It is challenging to judge the relation of the target entity pairs in this minimal supervision setting.

Word embedding, such as skip-gram (Mikolov et al., 2013a) and Glove (Pennington et al., 2014), are widely used in many natural language processing (NLP) tasks, and it was reported that word embeddings can capture the relational knowledge (Mikolov et al., 2013b). One intuitional method for relation induction task is using word embeddings to represent relations and induce relations based on vector translation or similarity (Vylomova et al., 2016; Drozd et al., 2016; Bouraoui et al., 2018; Vulić and Mrkšić, 2018; Camacho-Collados et al., 2019). However, the performance of these methods heavily depends on the pre-trained word embedding and these methods are rather noisy. According to the assumption that if two entities have a relationship in a known knowledge base, then all sentences that mention these two entities will express that relationship in some way (Mintz et al., 2009), many distant-supervised methods of relation extraction, such as PCNN(Zeng et al., 2015) and PCNN-BagATT (Ye and Ling, 2019) are proposed. Inspired by these methods, distant supervision might be another way to induce relation. To induce relation in the distant supervised way, we need a method to select proper sentences from corpus and extract relational knowledge from sentences. Luckily, many Pre-trained Language Models (PLMs), such as BERT(Devlin et al., 2019), GPT-2 (Radford et al., 2019) and XLNet(Yang et al., 2019), have been recently proposed and boost a great performance for many NLP tasks, such as question answering(Talmor et al., 2019; Feng et al., 2020), text summarization (Liu and Lapata, 2019; Lewis et al., 2020) and information extraction (Petroni et al., 2019; Alt et al., 2019). In order to better understand the PLMs, several works(Kim et al., 2020; Bouraoui et al., 2020; Ushio et al., 2021; Chen et al., 2021) have proven that PLMs can capture

084 syntactic and semantic knowledge. Bouraoui et al.
085 (2020) have explored the possibility of inducing
086 relation from BERT in a distant supervised way
087 and got a good result. To take the advantage that
088 BERT can capture context knowledge, they select
089 templates from corpus and fill entities in them to
090 let BERT predict the relation.

091 Existing methods are developed under the as-
092 sumption of sufficient seed entity pairs. However,
093 in practical scenarios, only a few entity pairs are
094 available for a particular relation. These methods
095 have difficulty in coping with the minimal super-
096 vision setting. The main reasons are: (1) Due to
097 the lack of labeled entity pairs, the model tends to
098 over-focus on the surface cues of the entity pairs
099 and ignores the contextual semantics. By simply
100 memorizing the seed entity pairs, it is difficult to
101 generalize the model to other entity pairs. (2) The
102 quality of templates is very important for relation
103 induction. When the seed entity pairs of a certain
104 relation are sparse, the number of candidate tem-
105 plates for this relation mined from the corpus will
106 be reduced.

107 Therefore, two major challenges should be ad-
108 dressed for the relation induction in the minimally-
109 supervised setting. (1) How to obtain a good gener-
110 alized relation induction model? (2) How to obtain
111 high-quality templates? so we propose a novel
112 approach called IST for minimally-supervised rela-
113 tion induction with Iteratively-Selected Templates
114 from PLM. Specifically, for the first challenge,
115 we use surface-agnostic features based on atten-
116 tion maps of BERT. Many works (Clark et al.,
117 2019; Kovaleva et al., 2019; Michel et al., 2019;
118 Wang et al., 2020) have revealed that the atten-
119 tion heads in BERT can capture much knowledge
120 and some attention heads are related to certain re-
121 lations, and some works use attention weights to
122 predict relations (Gu et al., 2021). For the second
123 challenge, we use Integrated Gradient (IG) (Sun-
124 dararajan et al., 2017) to score the templates and
125 iteratively select better templates. Intuitively, if a
126 sentence can well express the relational knowledge
127 between two entities, then the importance of these
128 two entities must be high in the sentence. On the
129 contrary, if a pair of entities do not play an impor-
130 tant role in a sentence, this sentence certainly does
131 not express the relationship between them. So IG
132 might be used to select high-quality sentences to
133 express relations.

134 We summarize our key contributions as follows:

- We propose a novel minimally-supervised re- 135
lation induction approach IST. To the best of 136
our knowledge, we are the first to address the 137
minimally-supervised relation induction task. 138
- In order to overcome the minimally- 139
supervised setting, we generate high-quality 140
pseudo-sentences by iteratively selecting 141
templates based on BERT and IG scores. 142
Moreover, we use attention maps to train a 143
more generalized model. 144
- We conduct extensive experiments on three 145
standard benchmark datasets, and our pro- 146
posed approach significantly outperforms the 147
state-of-the-art approaches. 148

2 Our Approach 149

150 In this section, we first formulate the minimally-
151 supervised relation induction task and give an
152 overview of our approach. We then describe the
153 details of each module in our approach.

2.1 Problem Formulation 154

155 Given a few seed entity pairs $P_r = \{(s_i, t_i)\}_{i=1}^N$
156 with a certain relation r , the task of relation in-
157 duction is to judge whether a new entity pair (s, t)
158 also has the relation r . In the minimally-supervised
159 setting, the number of the seed pairs is small for
160 each relation (in our experiments, no more than 5
161 per relation). To facilitate minimally-supervised re-
162 lation induction, we generate high-quality pseudo
163 sentences S_r for each relation r from a text corpus
164 C according to the seed entity pairs P_r with the
165 help of a pre-trained language model.

2.2 Overview 166

167 As illustrated in 1, our approach consists of
168 four main modules: template generation module,
169 pseudo sentence generation module, relation classi-
170 fier and template selection with IG.

171 In the template generation module, given
172 seed entity pairs, some proper templates could
173 be generated based on the mask-prediction re-
174 sults in BERT. For instance, considering the
175 seed entity pairs $P_r = \{(Germany, Berlin),$
176 $(France, Paris), (Japan, Tokyo)\}$, we can ob-
177 tain a sentence set S_r where each sentence men-
178 tions both entities of a pair in P_r . Taking a sentence
179 *The current capital of Japan is Tokyo.* as an
180 example, $\tau = (The\ current\ capital\ of\ _ is\ _)$ is
181 the generated template. Then, filling one entity into

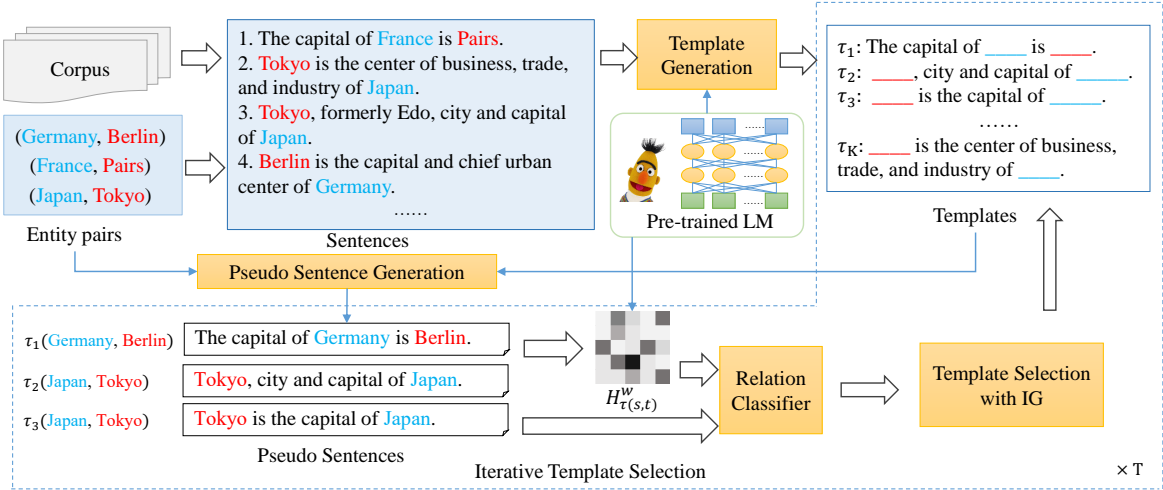


Figure 1: An overview of IST. First, we extract sentences that mention seed entity pairs as candidate sentences. Then, the template generation module uses BERT-prediction task to select proper templates from candidates. Templates and seed entity pairs are assembled to generate pseudo sentences to extract relational knowledge from BERT. The BERT attention weights between entities within the pseudo sentence are used as surface-agnostic features to better represent the relational knowledge. Then, the pseudo sentences and attention weights are combined to train a BERT-based relation classifier. Finally, we use the integrated gradient of entity pairs in pseudo sentences to evaluate the quality of templates and select better templates iteratively.

templates, the templates can be scored according to their ability to make BERT correctly predict another entity. This score is referred to as $score_{BERT}$.

After selecting proper templates based on $score_{BERT}$, we can generate pseudo sentences by assembling the templates and seed entity pairs. For example, a pseudo sentence $\tau(Germany, Berlin) = \text{The current capital of Germany is Berlin.}$ is generated by assembling $(Germany, Berlin)$ and τ . We generate both positive and negative sentences in this process.

For each pseudo-sentence, we extract surface-agnostic features based on attention weight maps of BERT and use them to train a relation classifier.

Finally, we use integrated gradient (IG) together with $score_{BERT}$ to evaluate the quality of templates again, so we can refine templates iteratively.

We will describe each module in detail in the following sections.

2.3 Template Generation

To induce relation from masked pre-trained language models such as BERT, we need templates for relations. First, many template-based relation extraction methods (Agichtein and Gravano, 2000; Ravichandran and Hovy, 2002) have proved that words near to s and t in corpus may represent a certain relation. To extract templates for relation r , we traverse Wiki Corpus to find k_i sentences

that contain both s_i and t_i ($i \leq N$), and the distance between s_i and t_i in sentence $D_{st} \leq d$. Then we mask s_i and t_i in these sentences to generate templates $\tau_{i,1}, \tau_{i,2}, \dots, \tau_{i,k}$. We can extract all candidate templates for r : $\{\tau_{1,1}, \tau_{1,2}, \dots, \tau_{i,j}, \dots, \tau_{N,k_N}\}$ ($i \leq N, j \leq k_i$), but not all of these templates are proper for inducing the relation r .

Then we need to select templates that are proper for BERT to induce relation r . Here we use BERT mask prediction as a template filter (Bouraoui et al., 2020). Specifically, for a template τ , insert s and t into τ respectively to get masked sentence $\tau(s, _)$ and $\tau(_, t)$. Then let BERT predict the masked token. If BERT can predict correctly, we consider the template τ is proper for relation r and $\tau(s, t)$ is “natural” for BERT.

$$score_{BERT}(\tau) = \sum_{i=0}^N (H(\tau(s_i, _)) + H(\tau(_, t_i))) \quad (1)$$

where $H(\tau(s_i, _))$ is 1 if the predicted token is t_i and 0 otherwise, and similar for $H(\tau(_, t_i)) = 1$.

By ranking templates with $score_{BERT}$, K proper templates $T_r = \{\tau_1, \tau_2, \dots, \tau_K\}$ are selected from candidate templates.

2.4 Pseudo Sentence Generation

In order to train a relation classifier, we assemble templates and seed entity pairs to generate labeled

235 pseudo sentences.

236 For positive sentences, we just assemble each
237 entity pair $(s, t) \in P_r$ with each template $\tau \in T_r$
238 to generate a sentence $\tau(s, t)$.

239 While for the negative sentences, follow-
240 ing(Vylomova et al., 2016), we have three strat-
241 egy for each pair $(s_i, t_i) \in P_r$. First, we ex-
242 change s and t as (t_i, s_i) (suppose r is not sym-
243 metrical). Second, we change one entity to another
244 entity in the same relation $:(s_i, t_j)$ or (s_j, t_i) ($i \neq$
245 $j, (s_i, t_i), (s_j, t_j) \in P_r$). Third, we change one
246 entity to an entity in other relations: (s_i, t_j) or
247 (s_j, t_i) ($i \neq j, (s_j, t_j) \in P_{r'}$).

248 2.5 Relation Classifier

249 Under the minimally-supervised setting, the model
250 should have good generalizability. We use surface-
251 agnostic features based on attention weights of
252 BERT to make model focus more on the relations
253 rather than the surface information of training data.

254 As Clark et al. (2019) has pointed out, some
255 heads of multi-head attention in BERT are related
256 to certain relations, and attention weights of cer-
257 tain heads can be used to extract certain relation
258 knowledge. Thus, for a proper template of relation
259 r , the attention weights between s and t of certain
260 heads related to r should be higher. But it is hard to
261 specify each head is related to what relations. Thus
262 we use attention weights of all heads as features to
263 induce relation knowledge.

264 Specifically, for a sentence $\tau(s, t)$, we calculate
265 the attention weights between s and t of all heads
266 as $\omega_{i,j,s \rightarrow t}$, where i denotes the i -th layer, j denotes
267 the j -th head in layer i and $s \rightarrow t$ denotes that this
268 is the attention s pays to t . We use the average
269 between $s \rightarrow t$ and $t \rightarrow s$ to express the attention
270 between them:

$$271 \omega_{i,j} = \frac{\omega_{i,j,s \rightarrow t} + \omega_{i,j,t \rightarrow s}}{2} \quad (2)$$

272 Then we construct attention weights embedding for
273 the sentence $\tau(s, t)$:

$$274 H_{\tau(s,t)}^{att} = \{\omega_{1,1}, \omega_{1,2} \dots \omega_{i,j} \dots \omega_{nl,nh}\} \quad (3)$$

275 where nl denotes the layer number, nh denotes
276 head number in a layer of BERT.

277 Besides $H_{\tau(s,t)}^{att}$, we also use BERT outputs to
278 represent the sentence $\tau(s, t)$. Specifically, we in-
279 put $\tau(s, t)$ into the BERT, and then use the output
280 vector of the [CLS] token as the feature $H_{\tau(s,t)}^{cls}$.
281 $H_{\tau(s,t)}^{cls}$ and $H_{\tau(s,t)}^{att}$ can compensate each other,

282 since $H_{\tau(s,t)}^{cls}$ can capture the information whether
283 $\tau(s, t)$ is ‘‘natural’’, and $H_{\tau(s,t)}^{att}$ contains the corre-
284 lation between (s, t) and the relation r . Thus, we
285 combine these two vectors through concatenation:

$$286 H_{\tau(s,t)} = H_{\tau(s,t)}^{cls} \oplus H_{\tau(s,t)}^{att} \quad (4)$$

287 Then, we feed $H_{\tau(s,t)}$ to a MLP classifier \mathcal{F}
288 and get the probability of (s, t) having relation r .
289 We use a cross-entropy loss to optimize \mathcal{F} . In
290 addition, we can also finetune BERT when training
291 the classifier.

292 2.6 Iterative Template Selection

293 BERT can rank templates by measuring whether a
294 sentence is natural. However, it can not capture the
295 different attribution of each token in a sentence for
296 expressing the relation.

297 Integrated Gradient is an attribution method pro-
298 posed in (Sundararajan et al., 2017).As Cui et al.
299 (2020) has described, the attribution score directly
300 reflects how much changing tokens will change the
301 model’s outputs. A higher attribution score repre-
302 sents more importance of tokens. In our relation
303 induction model, s and t obviously should be the
304 most important two tokens in sentences. Intuitively,
305 for a pseudo sentence $\tau(s, t)$, if the integrated gra-
306 dient value for s and t to the relation prediction is
307 higher, we are more confident that the relational
308 knowledge of (s, t) can be extracted well by the
309 model along with τ , so the template τ is much bet-
310 ter. Thus, we can use the integrated gradient of
311 (s, t) to the output of relation classifier to select
312 templates once again. Here, $\mathcal{F}(\tau, s, t)$ denotes the
313 relation classifier with $\tau(s, t)$ as the input.

314 According to Sundararajan et al. (2017), the in-
315 tegrated gradient value of s to $\mathcal{F}(\tau, s, t)$ is:

$$316 IG(\tau, s) = (s - s_0) \int_{x=0}^1 \frac{\partial \mathcal{F}(\tau, s_0 + \alpha(s - s_0), t)}{\partial s} d\alpha \quad (5)$$

317 where $\alpha \in [0, 1]$, and it can be approximated as:

$$318 IG(\tau, s) = (s - s_0) \sum_{i=1}^m \frac{1}{m} \times \frac{\partial \mathcal{F}(\tau, s_0 + \frac{i}{m}(s - s_0), t)}{\partial s} \quad (6)$$

319 where m is the number of approximate steps for
320 computing integrate gradient. For a template τ , we
321 calculate the average integrated gradient value for
322 all $(s, t) \in P_r$:

$$score_{IG}(\tau) = \sum_{(s,t) \in P_r} \frac{IG(\tau, s) + IG(\tau, t)}{2} \quad (7)$$

Then the templates are re-ranked according to the final score:

$$score = \alpha \cdot \frac{1}{rank_{BERT}} + (1 - \alpha) \frac{1}{rank_{IG}} \quad (8)$$

where $rank_{BERT}$ denotes the rank of templates according to $score_{BERT}$, $rank_{IG}$ denotes the rank of templates according to $score_{IG}$, and $\alpha \in [0, 1]$ is a coefficient to balance the two scores. Therefore, the templates could be selected iteratively for better relation induction.

2.7 Relation Induction

Given a new entity pair (x, y) , we fill them into templates τ_i , ($i \in K$) and use the classifier to predict $p_i(x, y)$, which denotes how much $\tau_i(x, y)$ is “natural”. Following Bouraoui et al. (2020), for all predictions from K templates $p_1(x, y), \dots, p_K(x, y)$, if $\max_i p_i(x, y) > 1 - \min_i p_i(x, y)$, then (x, y) is predicted to be positive.

3 Experiment Setup

3.1 Datasets

We conduct the experiments on three standard benchmark datasets in English: Google Analogy Test Set (GATS) (Mikolov et al., 2013a), Bigger Analogy Test Set (BATS) (Gladkova et al., 2016) and DiffVec (Vylomova et al., 2016).

GATS contains 5 semantic relations and 9 syntactic relations, and each consists of a varying number of entity pairs. While **BATS** contains 40 relations which are divided into 20 morphology relations and 20 semantic relations, each relation has 50 instances. **DiffVec** contains 36 relations with a various number of entity pairs. 10 of them are lexical or morphology relations and the remaining 26 are semantic relations.

3.2 Implementation Details

The relation induction task can be modeled as a binary classification problem for each relation. We first split the dataset into 50% of training data and 50% of test data. Then, under the minimally-supervised setting, for each relation r , we randomly select N entity pairs from training data as the seed entity pairs P_r . We extract candidate templates

from the English Wikipedia corpus¹ and $d = 15$. When generating K templates in T iterations, we initially select $K(T + 1)$ templates according to BERT score, and then iteratively filter out K improper templates in each iteration according to the score defined in Formula 8 until K templates are reserved at last for the final iteration. Notice that when $T = 0$, we only select K templates according to the $score_{BERT}$ without considering the $score_{IG}$. In our experiments, we use BERT-base², and set $N = 5, K = 20, T = 3, \alpha = 0.5$ by default.

We generate the same number of negative examples as positive examples for the training data and 3 times as many negative examples as positive examples for the test data.

For each relation, we repeat the experiments for 10 times and calculate the average result. The seed entity pairs used in each trial is randomly selected. The results of all metrics are calculated with micro-average.

4 Baselines

We compare our approach with three kinds of baselines.

The first kind is using the combination of pre-trained word embeddings to present relations. Specifically, following Vu and Shwartz (2018), we use $s \oplus t \oplus (s \odot t)$ to represent the relation between (s, t) and use a MLP classifier to make predictions. Here, the pre-trained word embeddings we used are Glove (Pennington et al., 2014)³ and SkipGram (Mikolov et al., 2013b)⁴. These two baselines are referred to as MLP_{sg} and MLP_{gl} respectively. We also use the Trans approach (Bouraoui et al., 2018) for relation induction by building subspaces for entities using word embeddings and modeling the relations with relative positions between subspaces.

The second kind is distant supervised methods. We use PCNN (Zeng et al., 2015) and PCNN-BagAtt (Ye and Ling, 2019) as two baselines. These distant supervised methods are proposed to solve the problem of noise in labeled data in relation extraction tasks. We also select the same number of sentences that mention entity pairs from the English Wikipedia corpus to construct training data

¹We used the dump of May 2021

²We used the BERT implementation available at <https://github.com/huggingface/transformers>

³<https://nlp.stanford.edu/projects/glove/>

⁴<https://code.google.com/archive/p/word2vec/>

	N=3									N=5								
	GATS			BATS			DiffVec			GATS			BATS			DiffVec		
	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1	Pre	Rec	F1
MLP _{sg}	41	54.4	39.6	40.3	45.8	40.6	40.1	49.8	41.5	43.3	56.8	45.3	43.5	47.1	43.8	43.5	51.3	43.9
MLP _{gl}	42.5	54.8	43.1	41.2	45.7	51.3	40.5	50.2	41.9	43.9	56.5	45.9	43.8	47.0	43.9	43.8	51.6	44.2
PCNN	58.6	52.5	56.1	52.1	45.8	47.7	57.3	51.5	53.9	60.1	56.2	58.3	53.4	45.8	50.3	59.0	52.7	55.4
PCNN_BagATT	63.9	56.2	59.5	56.4	45.9	48.6	61.5	52.8	55.8	65.3	58.6	60.1	57.8	51.0	80.8	63.5	53.4	57.1
BERT _{predict}	34.1	48.5	39.6	32.5	45.3	36.1	34.5	46.8	38.5	35.0	48.9	40.1	33.1	45.5	36.8	35.2	46.7	38.9
Trans	35.5	41.3	37.2	36.8	42.5	39.2	36.7	43.9	39.4	45.8	56.2	50.3	48.5	52.1	49.3	46.6	51.6	48.3
AutoPrompt	75.3	77.9	72.5	65.9	52.6	51.5	72.6	60.3	63.8	78.6	78.1	76.4	67.3	58.5	53.2	75.3	62.5	66.0
RI-BERT	79.3	80.5	75.8	70.1	53.0	53.2	77.4	65.8	68.3	80.7	80.1	79.5	70.1	55.7	55.9	79.5	67.2	70.4
IST	84.2	82.9	80.1	72.5	54.8	58.9	81.3	67.8	71.0	85.3	84.2	82.6	73.8	58.5	60.8	82.5	70.1	73.2

Table 1: Performance on three benchmarks when $N = 3$ and $N = 5$.

as in our approach. For an entity pair (s, t) , K sentences are used to predict its relation. If the average prediction score $\bar{S}_K \geq \theta$, (s, t) is predicted to have relation r . θ is a threshold and set to 0.7 in our experiments for its best performance.

The third kind is using the relational knowledge from PLMs, such as RI-BERT (Bouraoui et al., 2020), AutoPrompt (Shin et al., 2020)⁵ and BERT_{predict}.

RI-BERT induces relational knowledge from BERT, and our approach would degenerate to it when not using attention maps as the surface-agnostic features and not using $score_{IG}$ to refine the templates. We implement the method by ourselves since there is no open source.

AutoPrompt tries to elicit knowledge from PLM using automatically-constructed prompts. Here, We generate templates with AutoPrompt for each relation. Since there is only one template can be generated for each relation, we use a threshold-based method to determine whether a new entity pair (x, y) has a relation. When $p(x, y) > \delta$, the prediction would be positive. Here, $\delta = 0.8$ is the best threshold in our experiments.

BERT_{predict} is a simple baseline proposed by ourselves. After K templates are selected with $score_{BERT}$, we directly use BERT mask-prediction task to judge relation. Specifically, for an entity pair (s, t) and a template τ , if BERT can predict $\tau(s, _)$ or $\tau(_, t)$, the score of (s, t) will be increased by 1. The max score is $2K$, so if the score of $(s, t) \geq \epsilon \cdot 2K$, (s, t) is predicted to have the relation r . ϵ is a threshold and set to 0.7 for its best performance.

⁵<https://github.com/ucinlp/autoprompt>

5 Experimental Results

5.1 Main Results

The main experimental results on the three aforementioned benchmarks are shown in Table 1, which reports the micro-average of precision, recall and F1 of our approach **IST** and other state-of-the-art methods when $N = 3$ and $N = 5$.

From the table, there are several observations drawn from different aspects. (1) Our approach **IST** achieves the best performance against all other kinds of methods. (2) Pre-trained word embedding-based approaches such as MLP_{sg} and MLP_{gl} performance poorly, which proves that only few labeled entity pairs will degrade these approaches greatly. And Translation does not turn out well because of the lack of entities to construct representative subspaces. (3) The relational knowledge directly drawn from BERT also contains much noise according to the results of BERT_{predict}. (4) Traditional distant-supervised approaches which don't resort to PLM suffer from the noisy and sparse bag issues, although PCNN-BagATT uses intra-bag and inter-bag attention to handle sentence and bag-level noise, and get better performance, they are still not suitable for the minimally-supervised relation induction task. (5) AutoPrompt and RI-BERT use proper prompts or templates from BERT, so they can obtain a better performance. However, they did not consider the generalization problem in the minimally-supervised setting. In addition, they ignored the contribution of each token in a sentence for expressing the relation, especially for the entity pairs, but only considered whether a sentence is natural or not according to BERT. (6) More labeled entity pairs can achieve better performance by comparing the results of $N = 3$ and $N = 5$. This

phenomenon is reflected by all methods in both three datasets.

5.2 Ablation Study and Analysis

Performance of Different Relations To further explore the performance of different relations, we show the detailed results of each relation in GATS in Table 2.

From the table, we can see that our approach achieves better performance for both semantic and morphology relations. Moreover, the iteratively template selection can bring a significant improvement, especially for semantic relations. As to morphological relations, the improvement is not so evident. This is because the entities in morphological relations are always adverbs or adjectives to which little attention is paid, so $H_{\tau(s,t)}^{att}$ plays a limited role.

	GATS	RI-BERT	T=0	T=1	T=2
Semantic	currency	56.7	58.8	58.6	59.5
	family	76.9	78.8	78.4	79.9
	capital-common	88.4	87.3	85.7	91.6
	city-in-state	68.2	71.0	73.1	75.2
	capital-world	77.3	76.8	78.0	78.2
	Average	73.5	74.5	74.7	76.9
Morphology	adj-to-adv	39.1	38.8	42.3	44.8
	opposite	55.3	59.7	54.0	56.6
	comparative	90.9	87.5	88.2	89.0
	superlative	78.1	79.8	80.6	77.7
	presen-participle	98.4	96.2	98.1	98.9
	nationality-adj	91.5	92.4	91.7	92.1
	past-tense	96.9	97.8	97.2	97.0
	plural	93.8	91.6	96.6	95.8
	plural-verb	100	99.0	99.7	99.7
	Average	82.6	82.6	83.2	83.5

Table 2: Detailed experimental results (F1) for each relation on GATS.

Performance of attention weights and IG To investigate the effectiveness of BERT attention weights and IG, we compare the performance of several variants of our approach on GATS.

To reduce the effect of BERT attention weights, the representation of sentence $\tau(s, t)$ is simplified from $H_{\tau(s,t)}^{cls} \oplus H_{\tau(s,t)}^{att}$ to $H_{\tau(s,t)}^{cls}$. In addition, without IG, there would be no iterative template selection procedure. The results are shown in Table 3, and the performance drops in all variants, which proves that both attention maps and integrated gradient are useful in our approach.

Different Number of Templates To analyze the impacts of the number of templates (K), we conduct experiments with different numbers of templates, and the results are shown in Table 4. From

GATS	T=0	T=1	T=2	T=3
IST	79.7	80.2	81.1	82.6
w/o att	79.5	79.7	80.6	80.9
w/o IG	78.4	78.4	78.4	78.4

Table 3: The F1 scores of IST and other variants on GATS with different iterations.

the table, we find that more templates can bring better performance in all iterations. However, if K is too large, the time consumption will be greater and some unsuitable templates will be retained, leading to worse results.

	K=5	K=10	K=20
T=0	72.5	75.3	79.7
T=1	73.8	78.5	80.2
T=2	75.6	78.9	81.1

Table 4: F1-score with different number of templates ($K = 5, 10, 20$) and different iterations ($T = 0, 1, 2$) on GATS.

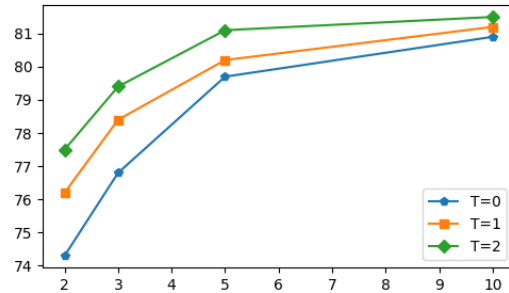


Figure 2: F1 scores with different numbers of seed entity pairs ($N = \{2, 3, 5, 10\}$) of our approach on GATS

Different Number of Seeds We evaluate our approach with different numbers of seed entity pairs (N), and the results are shown in Figure 2. From the figure, we can see that F1 score increases gradually until convergence for all iterations. Our approach already achieves a satisfactory result when $N = 5$.

Effect of Balance Coefficient The parameter $\alpha \in [0, 1]$ is a balance coefficient between $score_{BERT}$ and $score_{IG}$ for template scoring. Larger α will consider $score_{IG}$ more in the scoring. We conduct the experiments with different α on GATS, and the results are shown in Figure 3. From the figure, we find that our approach achieves the best performance when $\alpha = 0.5$.

T=0	T=3
The Government of _ denoted 300 million _ to finance the school’s construction in 1975.	The _ (, plural: / ,) is the currency of _ .
Currently, _ uses the _ as its national currency.	This was one of the reasons for naming the current currency of the Republic of _ the _.
Following the introduction of the euro, the _ was linked to the euro, until January 1, 2015, when _ officially adopted the euro as its currency.	The _ (; sign: ; code: KHR) is the currency of _.
AutoPrompt: _ cial largest greenwich _.	

Table 5: Case study for relation *currency*, where top 3 templates are exhibited with different approaches.

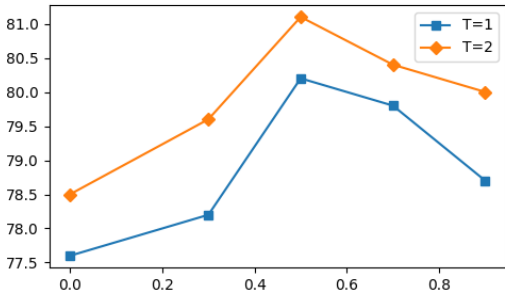


Figure 3: F1 scores with $\alpha=\{0, 0.3, 0.5, 0.7, 0.9\}$ of our approach on GATS

Case Study Table 5 compares the selected templates of relation *currency* between $T = 0$ and $T = 3$. From human’s intuition, we find that comparing to ($T = 3$), the templates filtered out only with $score_{BERT}(T = 0)$ are more ambiguous that they might indicate a co-occurrence relationship rather than relation *currency*. For example, for the first template “The Government of _ denoted 300 million _ to finance the school’s construction in 1975.”, it is natural for the government of a country to denote their own currency or just use dollar to evaluate how much they have denoted. So $\tau(s, dollar)$ is natural when s denotes any country. This is due to the way of selecting templates that only requires the templates is proper for all $(s, t) \in P_r$ without explicitly declaring what the relation is. In fact, the model can distinguish *co-occurrence* and *currency* only after the BERT is fine-tuned with negative examples. As to the template generated by AutoPrompt, it is a combination of some tokens rather than a human-readable sentence. Although AutoPrompt got good results on some tasks(Shin et al., 2020), the template is totally not interpretable from human’s perspective.

6 Related Work

6.1 Relation Induction

Relation induction was first proposed in (Vylomova et al., 2016). They used the vector difference between two entities to represent the relation between them. More researches on the relation induction with word embeddings were proposed in(Drozd et al., 2016; Bouraoui et al., 2018; Vu and Shwartz, 2018). They pointed out that the difference is not the best way to express the relationship and proposed more complicated methods to better extract relational knowledge between word embeddings.

6.2 Knowledge Induction from BERT

BERT was proven to be able to capture relational knowledge(Kim et al., 2020; Bouraoui et al., 2020; Ushio et al., 2021; Chen et al., 2021).Inspired by this, some works tried to use BERT on the relation induction task (Shin et al., 2020; Bouraoui et al., 2020; Jiang et al., 2020).The key point of these methods is to fill entities in the proper templates.

Recently, many efforts focus on the generation of templates. Jiang et al. (2020) proposed a template generation strategy based on paraphrasing aiming to improve lexical diversity while remaining relatively faithful to the original prompt. Shin et al. (2020) proposed AutoPrompt method to generate templates, or as they called , prompts, from nothing instead of from corpus. They automated create prompts based on gradient-guided search.

7 Conclusion

In this paper, we propose a novel minimally-supervised relation induction approach. Our proposed approach can iteratively select proper templates using $score_{IG}$ and $socre_{IG}$, and obtain a good generalized ability with surface-agnostic features based on attention maps of BERT. Experiments illustrate that our approach achieves state-of-the-art performance on three standard benchmarks.

597
598
599
600
601
602

603
604
605
606
607
608
609

610
611
612
613

614
615
616
617
618
619

620
621
622
623
624
625

626
627
628
629
630
631
632

633
634
635
636
637
638
639

640
641
642

643
644
645
646
647
648
649
650
651

References

Eugene Agichtein and Luis Gravano. 2000. [Snowball: Extracting relations from large plain-text collections](#). In *Proceedings of the Fifth ACM Conference on Digital Libraries*, DL '00, page 85–94, New York, NY, USA. Association for Computing Machinery.

Christoph Alt, Marc Hübner, and Leonhard Hennig. 2019. [Fine-tuning pre-trained transformer language models to distantly supervised relation extraction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1388–1398, Florence, Italy. Association for Computational Linguistics.

Zied Bouraoui, Jose Camacho-Collados, and Steven Schockaert. 2020. [Inducing relational knowledge from bert](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7456–7463.

Zied Bouraoui, Shoaib Jameel, and Steven Schockaert. 2018. [Relation induction in word embeddings revisited](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1627–1637, Santa Fe, New Mexico, USA. Association for Computational Linguistics.

Jose Camacho-Collados, Luis Espinosa Anke, and Steven Schockaert. 2019. [Relational word embeddings](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3286–3296, Florence, Italy. Association for Computational Linguistics.

Catherine Chen, Kevin Lin, and Dan Klein. 2021. [Constructing taxonomies from pretrained language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4687–4700, Online. Association for Computational Linguistics.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.

Leyang Cui, Sijie Cheng, Yu Wu, and Yue Zhang. 2020. [Does bert solve commonsense task via commonsense knowledge?](#) *ArXiv*, abs/2008.03945.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Aleksandr Drozd, Anna Gladkova, and Satoshi Matsuoka. 2016. [Word embeddings, analogies, and machine learning: Beyond king - man + woman = queen](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3519–3530, Osaka, Japan. The COLING 2016 Organizing Committee.

Yanlin Feng, Xinyue Chen, Bill Yuchen Lin, Peifeng Wang, Jun Yan, and Xiang Ren. 2020. [Scalable multi-hop relational reasoning for knowledge-aware question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1295–1309, Online. Association for Computational Linguistics.

Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. 2016. [Analogy-based detection of morphological and semantic relations with word embeddings: what works and what doesn’t](#). In *Proceedings of the NAACL Student Research Workshop*, pages 8–15, San Diego, California. Association for Computational Linguistics.

Xiaotao Gu, Zihan Wang, Zhenyu Bi, Yu Meng, Liyuan Liu, Jiawei Han, and Jingbo Shang. 2021. [Ucphrase: Unsupervised context-aware quality phrase tagging](#). *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*.

Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. [How can we know what language models know?](#) *Transactions of the Association for Computational Linguistics*, 8:423–438.

Taeuk Kim, Jihun Choi, Daniel Edmiston, and Sang goo Lee. 2020. [Are pre-trained language models aware of phrases? simple but strong baselines for grammar induction](#). *ArXiv*, abs/2002.00737.

Olga Kovaleva, Alexey Romanov, Anna Rogers, and Anna Rumshisky. 2019. [Revealing the dark secrets of BERT](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4365–4374, Hong Kong, China. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Yang Liu and Mirella Lapata. 2019. [Text summarization with pretrained encoders](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.

710	Paul Michel, Omer Levy, and Graham Neubig. 2019.	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and	764
711	Are sixteen heads really better than one? In <i>NeurIPS</i> .	Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In <i>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)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.	765
712	Tomas Mikolov, Kai Chen, Gregory S. Corrado, and	Asahi Ushio, Luis Espinosa Anke, Steven Schockaert,	766
713	Jeffrey Dean. 2013a. Efficient estimation of word	and Jose Camacho-Collados. 2021. BERT is to NLP what AlexNet is to CV: Can pre-trained language models identify analogies? In <i>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)</i> , pages 3609–3624, Online. Association for Computational Linguistics.	767
714	representations in vector space. In <i>ICLR</i> .		768
715	Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig.	Tu Vu and Vered Shwartz. 2018. Integrating multiplicative features into supervised distributional methods for lexical entailment. In <i>Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics</i> , pages 160–166, New Orleans, Louisiana. Association for Computational Linguistics.	769
716	2013b. Linguistic regularities in continuous space		770
717	word representations. In <i>Proceedings of the 2013</i>	Ivan Vulić and Nikola Mrkšić. 2018. Specialising word	771
718	<i>Conference of the North American Chapter of the</i>	vectors for lexical entailment. In <i>Proceedings of</i>	772
719	<i>Association for Computational Linguistics: Human</i>	<i>the 2018 Conference of the North American Chapter</i>	773
720	<i>Language Technologies</i> , pages 746–751, Atlanta,	<i>of the Association for Computational Linguistics:</i>	774
721	Georgia. Association for Computational Linguistics.	<i>Human Language Technologies, Volume 1 (Long</i>	775
722	Mike Mintz, Steven Bills, Rion Snow, and Daniel Ju-	<i>papers)</i> , pages 1134–1145, New Orleans, Louisiana.	776
723	rafsky. 2009. Distant supervision for relation ex-	Association for Computational Linguistics.	777
724	traction without labeled data. In <i>Proceedings of the</i>		778
725	<i>Joint Conference of the 47th Annual Meeting of the</i>		779
726	<i>ACL and the 4th International Joint Conference on</i>		780
727	<i>Natural Language Processing of the AFNLP</i> , pages		781
728	1003–1011, Suntec, Singapore. Association for Com-		
729	putational Linguistics.		
730	Jeffrey Pennington, Richard Socher, and Christopher		
731	Manning. 2014. GloVe: Global vectors for word		
732	representation. In <i>Proceedings of the 2014 Confer-</i>		
733	<i>ence on Empirical Methods in Natural Language Pro-</i>		
734	<i>cessing (EMNLP)</i> , pages 1532–1543, Doha, Qatar.		
735	Association for Computational Linguistics.		
736	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,		
737	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and		
738	Alexander Miller. 2019. Language models as knowl-		
739	edge bases? In <i>Proceedings of the 2019 Confer-</i>		
740	<i>ence on Empirical Methods in Natural Language Pro-</i>		
741	<i>cessing and the 9th International Joint Conference</i>		
742	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,		
743	pages 2463–2473, Hong Kong, China. Association		
744	for Computational Linguistics.		
745	Alec Radford, Jeff Wu, Rewon Child, David Luan,		
746	Dario Amodei, and Ilya Sutskever. 2019. Language		
747	models are unsupervised multitask learners.		
748	Deepak Ravichandran and Eduard Hovy. 2002. Learn-		
749	ing surface text patterns for a question answering		
750	system. In <i>Proceedings of the 40th Annual Meet-</i>		
751	<i>ing of the Association for Computational Linguistics</i> ,		
752	pages 41–47, Philadelphia, Pennsylvania, USA. As-		
753	sociation for Computational Linguistics.		
754	Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric		
755	Wallace, and Sameer Singh. 2020. AutoPrompt: Elic-		
756	iting Knowledge from Language Models with Auto-		
757	matically Generated Prompts. In <i>Proceedings of the</i>		
758	<i>2020 Conference on Empirical Methods in Natural</i>		
759	<i>Language Processing (EMNLP)</i> , pages 4222–4235,		
760	Online. Association for Computational Linguistics.		
761	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.		
762	Axiomatic attribution for deep networks. <i>ArXiv</i> ,		
763	abs/1703.01365.		

821 *in Natural Language Processing*, pages 1753–1762,
822 Lisbon, Portugal. Association for Computational Lin-
823 guistics.

824 **A Example Appendix**

825 This is an appendix.