LREC-COLING 2024



Discriminative Language Model as Semantic Consistency Scorer for Prompt-based Few-Shot Text Classification

Zhipeng XIE, Yahe LI

School of Computer Science, Fudan University, China xiezp@fudan.edu.cn

Outline

Background

- Method
- Experimental Results

• Future Work

Pretraining-Finetuning Paradigm

Pretraining-Finetuning paradigm - de facto standard for NLU and NLG tasks

Pretrained Language Models:

Autoregressive Models (ARM) – Predict next token based on previous ones

✓ Effective for NLG tasks

Masked Language Model (MLM) – Reconstruct masked tokens based on their bidirectional surrounding contexts

✓ Good at NLU tasks but usually not applicable to NLG

Finetuning PLM for Few Shot Text Classification

- Few Shot Text Classification: The amount of annotated examples is limited.
- Two types of finetuning techniques:
 - Conventional Finetuning
 - Prompt-based Finetuning
 - Masked Language Model
 - Discriminative Language Model

Conventional Finetuning



Conventional finetuning works well with abundant training examples, but will be cornered in the few shot scenario.

Prompt-based Finetuning

Prompt-based prediction was originally developed by the GPT series for zero-shot predicitons

- (Radford et al., 2018, 1029; Brown et al, 2020)
- PET method studied prompt-based prediction for finetuning
 - (Schick and Schutze, 2021a,b)
- LM-BFF method automated the process of prompt generation
 - (Gao et al., 2021)

Prompt-based Finetuning for Text Classification (3) The model makes prediction according to the probabilities of filling the [MASK] token with the label words



ELECTRA – A Discriminative LM (Clark et al., 2020)



- The pretraining task: replaced token prediction.
 - The **generator** is trained to perform MLM task
 - The **discriminator** is trained to distinguish "real" tokens from "replaced" tokens
- The discriminator in **ELECTRA** consists of two modules:
 - An encoder: maps a sequence of input tokens $x = [x_1, ..., x_n]$ into a sequence of contextualized vector representations $[\mathbf{h}_1, ..., \mathbf{h}_n]$
 - A discriminative head: predicts whether each token is a "real" or "replaced" token.

Prompting DLM for Text Classification DPT (Yao et al., 2022)

• The DPT method converts each input text *x* into the following template of discriminative prompt:

[CLS] x Class: $v(l_1), v(l_2), \dots, v(l_n)$. [SEP]

Input example: The restaurant has excellent food.

Discriminative Prompt: [CLS] The restaurant has excellent food. Class: great, terrible. [SEP]

The prompt in this style is **not natural** and is **not compatible** with the training data used for pretraining.

Prompting DLM for Text Classification PromptELECTRA (Xia et al., 2022)

- **PromptELECTRA** generates one discriminative prompt for each possible class label.
- DLM head is used to output the label word that has the highest probability of being original token in its corresponding prompt.

The prompts in this style is much more natural. But this method simply rely on the evidence from the label words, other possible evidences get ignored

Outline

- Introduction
- Method
- Experimental Results
- Future Work

Motivation <

Input sentence: x = "The restaurant has excellent food."



Method

- The **DLM-SCS** method stands for **Discriminative Language** <u>Model as Semantic Consistency Scorer</u>
- Basic Idea: reformulate text classification as a task of semantic consistency scoring. The class label with the highest semantic consistency score is predicted:

$$\hat{l} = \operatorname*{arg\,max}_{l \in \mathcal{L}} SC(\tilde{\boldsymbol{x}}^l)$$

- Technique: How to calculate the semantic consistency of each discriminative prompt?
 - The weighted average of semantic consistency scores of multiple parts in the prompt.

Discriminative Prompts ____ Sentence Classification

- Given an input example of sentence $x_{in} = x^{(1)}$
- The discriminative prompt for each label $l \in \mathcal{L}$ is generated as: $\widetilde{x}^{l} = [\text{CLS}] \ x^{(1)} It \text{ is } v(l) . [\text{SEP}]$
- Each discriminative prompt \widetilde{x}^l has two parts:
 - The sentence $x^{(1)}$
 - The label word v(l)

Discriminative Prompts Sentence Pair Classification

- Given an input example of sentence pair $x_{in} = (x^{(1)}, x^{(2)})$
- The discriminative prompt for each label $l \in \mathcal{L}$ is generated as: $\widetilde{x}^{l} = [\text{CLS}] \ x^{(1)} ? v(l), x^{(2)} [\text{SEP}]$
- Each discriminative prompt \widetilde{x}^l has three parts:
 - The first sentence $\pmb{x^{(1)}}$
 - The second sentence $x^{(2)}$
 - The label word v(l)

Semantic Consistency Score: Discriminative Prompt

Semantic consistency score of a discriminative prompt



Semantic Consistency Score <a>
 Token Subsequence — Uniform Weights

Let s be a token subsequence in the discriminative prompt \tilde{x}^l , its semantic inconsistency is measured as:

$$sc(\boldsymbol{s}, \tilde{\boldsymbol{x}}^{\boldsymbol{l}}) = \frac{\exp\left(-\frac{1}{|\boldsymbol{s}|} \sum_{x \in \boldsymbol{s}} \boldsymbol{w}^{\top} \boldsymbol{h}_{x}^{\boldsymbol{l}}\right)}{\sum_{l' \in \mathcal{L}} \exp\left(-\frac{1}{|\boldsymbol{s}|} \sum_{x \in \boldsymbol{s}} \boldsymbol{w}^{\top} \boldsymbol{h}_{x}^{l'}\right)}$$

Semantic Consistency Score <a>Token Subsequence – IDF Weights

Different tokens should be of different importance with respect to the semenatic consistency

using IDF value as the weights:

$$sc(\boldsymbol{s}, \tilde{\boldsymbol{x}}^{l}) = \frac{\exp\left(-\frac{\sum_{x \in \boldsymbol{s}} \mathsf{idf}(x) \mathbf{w}^{\top} \mathbf{h}_{x}^{l}}{\sum_{x \in \boldsymbol{s}} \mathsf{idf}(x)}\right)}{\sum_{l' \in \mathcal{L}} \exp\left(-\frac{\sum_{x \in \boldsymbol{s}} \mathsf{idf}(x) \mathbf{w}^{\top} \mathbf{h}_{x}^{l'}}{\sum_{x \in \boldsymbol{s}} \mathsf{idf}(x)}\right)}$$

DLM-SCS: Schematic Illustration



(a) Prompting discriminative language model as semantic consistency scorer

Outline

- Background
- Method
- Experimental Results
- Future Work

Experimental Setup Datasets

- Datasets:
 - 4 Sentence Classification Datasets: SST-2, SST-5, MR and CR
 - 6 Sentence-Pair Classification Datasets: SNLI, MNLI, QNLI, RTE, MRPC and QQP

- The number of training examples:
 - K = 16 (by default) training examples per class
 - The total number of training examples is $K \times |\mathcal{L}|$
 - Development set for model selection and hyperparameter tuning is of the same size as the training set

Experimental Setup Templates and Label Words

• We adopt the templates and label words from (Gao et al., 2021):

Task	Template	Label words			
SNLI	$<\!\!S_1\!\!>? v(l), <\!\!S_2\!\!>$	Yes/No/Maybe			
MNLI	$<\!S_1\!>? v(l), <\!S_2\!>$	Yes/No/Maybe			
QNLI	$<\!S_1\!>? v(l), <\!S_2\!>$	Yes/No			
RTE	$<\!S_1\!>? v(l), <\!S_2\!>$	Yes/No			
MRPC	$<\!S_1\!>? v(l), <\!S_2\!>$	Yes/No			
QQP	$<\!S_1\!>. v(l), <\!S_2\!>$	Yes/No			
SST-2	$< S_1 >$ It is $v(l)$.	terrible/great			
SST-5	$\sim S_{1} > t is w(1)$	terrible/bad/			
	$\nabla_1 \ge 1$ is $v(i)$.	okay/good/great			
MR	$< S_1 >$ It is $v(l)$.	terrible/great			
CR	$< S_1 >$ It is $v(l)$.	terrible/great			

Experimental Setup The Competitors

- Fine-tuning: the traditional fine-tuning of Roberta-Large
- LM-BFF (man): few-shot finetuning with manual prompts
- LM-BFF (auto): few-shot finetuning with automatically searched templates.

- Two prompt methods for discriminative language model:
 - DPT
 - PromptELECTRA

	Model	SNLI	MNLI	QNLI	RTE	MRPC	QQP	SST-2	SST-5	MR	CR	
	Woder	(acc)	(acc)	(acc)	(acc)	(F1)	(F1)	(acc)	(acc)	(acc)	(acc)	
	Fine-tuning	48.4	45.8	60.2	54.4	76.6	60.7	81.4	43.9	76.9	75.8	
		(4.8)	(6.4)	(6.5)	(3.9)	(2.5)	(4.3)	(3.8)	(2.0)	(5.9)	(3.2)	
	LM-BFF (man)	77.2	68.3	64.5	69.1	74.5	<mark>65.5</mark>	92.7	47.4	87.0	90.3	
		(3.7)	(2.3)	(4.2)	(3.6)	(5.3)	<mark>(5.3)</mark>	(0.9)	(2.5)	(1.2)	(1.0)	
	+demonstrations	79.7	70.7	69.2	68.7	77.8	69.8	92.6	50.6	86.6	90.2	
Mean performance of the runner-up: 72.7		<mark>(1.5)</mark>	(1.3)	(1.9)	(2.3)	(2.0)	(1.8)	(0.5)	(1.4)	(2.2)	(1.2)	
	LM-BFF (auto)	77.1	68.3	68.3	73.9	76.2	67.0	92.3	49.2	85.5	89.0	
		<u>(2</u> .1)	<u>(2.5)</u>	<u>(7.4)</u>	<u>(2.2)</u>	<u>(2.3)</u>	<u>(3.0)</u>	(1.0)	(1.6)	(2.8)	(1.4)	Moon
	+demonstrations	77.5	70.0	68.5	<mark>71.1</mark>	78.1	67.7	93.0	49.5	87.7	91.0	
		(3.5)	(3.6)	(5.4)	(5.3)	(3.4)	(5.8)	(0.6)	(1.7)	(1.4)	(0.9)	performance
	DART	75.8	67.5	66.7	68.7	78.3	67.8	93.5	49.6	88.2	91.8	of the
		(1.6)	(2.6)	(3.7)	(1.3)	(4.5)	(3.2)	(0.5)	(0.9)	(1.0)	(0.5)	runner-up:
	DPT	47.4	39.0	54.6	50.2	76.4	56.1	92.6	44.0	89.5	91.2	81.0
	2	(7.7)	(1.8)	(5.4)	(2.8)	(6.1)	(1.1)	(1.3)	(<u>3.8</u>)	(2.1)	(1.6)	1
	PromptELECTRA	79.1	65.8	70.9	68.2	73.5	63.1	93.1	51.4	89.4	90.2	
iviean		(3.4)	(2.5)	(2.1)	(2.8)	(4.6)	<mark>(3.3)</mark>	(1.0)	(2.2)	(1.6)	(1.4)	
Performance	DLM-SCS (ours)	82.2	71.0	77.0	75.0	78.3	72.2	93.6	51.5	90.2	91.0	
of Our		(1.5)	(2.0)	(2.4)	(2.9)	(3.1)	(1.4)	(0.6)	(2.0)	(0.7)	(1.4)	
method: 76.0												

Mean Performance of Our method: 81.6

Ablation Analysis

Two main techniques in DLM-SCS:

1) Integrating the evidences from multiple components (or parts) of the prompt

2) weighting the tokens in each prompt part with IDF values

Model	SNLI	MNLI	QNLI	RTE	MRPC	QQP	SST-2	SST-5	MR	CR
(full) DLM-SCS	82.2	71.0	77.0	75.0	78.3	72.2	93.6	51.5	90.2	91.0
-w.o. token weight	78.2	70.0	73.4	73.6	76.9	69.6	93.0	48.8	90.3	90.3
-only label word	76.5	64.6	69.0	71.8	74.8	64.2	93.7	51.1	88.8	90.4



Varying the size of training examples



Outline

- Background
- Method
- Experimental Results
- Future Work

Two Future Directions

- How to integrate our discriminative framework with auotomatic prompt generation and differentiable prompting
- How to combine our discriminative framework with generative PLMs?

Thanks for Listening! Any Questions?