

Decoupling Knowledge from Memorization: Retrieval-augmented Prompt Learning



Xiang Chen ^{1,2}, Lei Li ^{1,2}, Ningyu Zhang ^{1,2} *, Xiaozhuan Liang ^{1,2}, Shumin Deng ^{1,2}

Chuanqi Tan ³, Fei Huang ³, Luo Si ³, Huajun Chen ^{1,2}*

¹Zhejiang University & AZFT Joint Lab for Knowledge Engine,

²Hangzhou Innovation Center, Zhejiang University, ³Alibaba Group

Introduction

• Limitations of Prompt Learning:



different Prompt learning with PLMs usually generalizes unstably in an extremely lowresource setting or emerging domains. One potential reason is that, it is non-trivial for parametric models to learn rare or hard patterns well with rote memorization, thus, resulting in inefficient generalizable performance.



• Decoupling knowledge from memorization :

with the motivation of decoupling knowledge from memorization to help the model strike a balance between generalization and memorization, we constructs an open-book knowledgestore from training instances and implements a retrieval mechanism during the process of input, training and inference, thus equipping the model with the ability to retrieve related contexts from the training corpus as cues for enhancement.

Retrieval-augmented Prompt Learning

Open-book Knowledge-store

Given the *i*-th example (c_i , y_i) in the training data C, we obtain the key-value pair $(\mathbf{h}_{\hat{c}_i}, v_i)$, in which $\hat{\mathbf{c}}_i = \mathcal{T}(\mathbf{c}_i), \mathbf{h}_{\hat{c}_i} \in \mathbb{R}^d$ is the embedding of the [MASK] token in the last layer of the PLM, and $v_i = f(y_i)$ denotes the label word of the *i*th example.

 $(\mathcal{K},\mathcal{V}) = ig\{ oldsymbol{h}_{\hat{oldsymbol{c}}_i},v_iig) \mid (oldsymbol{c}_i,y_i) \in \mathcal{C}ig\}$



Retrieval of Neural Demonstration

We intuitively aggregate the *m* neighbor vectors for each class according to their similarity and incorporate the demonstration into the input representation of \hat{x} after the word embedding layer of the \mathcal{M} as follows:

$$\mathcal{I} = e(\hat{m{x}}) \oplus \left[\sum_{i \in [1:m]} lpha_i^{(1)} m{h}_{\hat{c}_i}^{(1)}, e\Big(v^{(1)}\Big)
ight] \oplus \ldots \oplus \left[\sum_{i \in [1:m]} lpha_i^{(L)} m{h}_{\hat{c}_i}^{(L)}, e\Big(v^{(L)}\Big)
ight]; lpha_i^{(l)} = rac{e^{m{h}_{\hat{q}} \cdot m{h}_{\hat{c}_i}^{(l)}}}{\sum_{i \in [1:m]} e^{m{h}_{\hat{q}} \cdot m{h}_{\hat{c}_i}^{(l)}}}$$

Retrieve kNN for Guiding Training

a. Retrieval-augmented prompt learning

b. Creation and refresh of open-book knowledge-store

Figure 2: Overview of RETROPROMPT. Note that $e(\cdot)$ denotes word embedding function in the PLM \mathcal{M} , while "M", "t" and "g" in $e(\cdot)$ specifically refers to "[MASK]", "terrible" and "great".

Our intuition is to differentiate between easy and hard examples according to the prediction of kNN.

$$egin{aligned} P_{k ext{NN}}(y \mid oldsymbol{q}_t) \propto \sum_{(oldsymbol{c}_i, y_i) \in \mathcal{N}} oldsymbol{1}_{y=y_i} \expig(dig(oldsymbol{h}_{oldsymbol{\hat{q}}_t}, oldsymbol{h}_{oldsymbol{\hat{c}}_i}ig)ig) \ F(p_{k ext{NN}}) &= -\log(p_{k ext{NN}}), \quad \mathcal{L} = (1+eta F(p_{k ext{NN}}))\mathcal{L}_{CE} \end{aligned}$$

kNN based probability for Cloze-style Prediction

we reformulate the $P(y | q_t)$ by interpolating the P_{kNN} with the already-trained base PLM's MLM prediction $P_{\mathcal{M}}$ using parameter λ to produce the final probability of the label:

$$P(y \mid oldsymbol{q}_t) = \lambda P_{k\mathrm{NN}}(y \mid oldsymbol{q}_t) + (1 - \lambda)g(P_{\mathcal{M}}([\mathrm{MASK}] = v \mid \mathcal{T}(oldsymbol{q}_t)))$$

Experiments

Few-shot/Zero-shot Results

St.	Model	Single Sentence			Sentence Pair				Information Extraction			
		SST-2 (acc)	MR (acc)	CR (acc)	MNLI (acc)	QNLI (acc)	QQP (F1)	Model	FewN (acc)	SemEval (acc)	TACRED (F1)	Avg.
16	FT LM-BFF (man) LM-BFF (D-demo) KPT †	81.4 (3.8) 91.6 (1.2) 91.8 (1.2) 90.3 (1.6)	76.9 (5.9) 87.0 (2.0) 86.6 (1.8) 86.8 (1.8)	75.8 (3.2) 90.3 (1.6) 90.2 (1.4) 88.8 (3.7)	45.8 (6.4) 64.3 (2.5) 64.8 (2.3) 61.4 (2.1)	60.2 (6.5) 64.6 (5.4) 69.2 (5.4) 61.5 (2.8)	60.7 (4.3) 65.4 (5.3) 68.2 (3.2) 71.6 (2.7)	FT KnPr KnPr (D-demo) KPT †	52.7 (2.2) 65.3 (1.1) 65.9 (1.5)	66.1 (1.2) 80.9 (2.5) 78.8 (2.1)	25.8 (2.8) 33.2 (2.0) 32.8 (1.7)	60.6 71.4 72.2* 70.9
	Ours	93.9 (0.4)	88.0 (0.8)	91.9 (0.7)	71.1 (1.8)	71.6 (1.8)	74.0 (2.0)	Ours	67.3 (0.9)	81.5 (1.3)	40.7 (0.7)	75.6
4	FT LM-BFF (man) LM-BFF (D-demo) KPT †	60.2 (2.8) 90.7 (0.8) 90.2 (1.5) 88.2 (5.7)	57.6 (1.4) 85.2 (2.8) 85.5 (2.1) 83.4 (1.5)	66.4 (5.5) 89.9 (1.8) 89.7 (0.6) 87.2 (2.5)	35.0 (0.3) 51.0 (2.5) 56.1 (1.0) 53.7 (2.7)	54.2 (3.9) 61.1 (6.1) 61.7 (7.6) 59.2 (2.8)	52.8 (4.7) 48.0 (4.9) 63.2 (5.6) 54.9 (7.9)	FT KnPr KnPr (D-demo) KPT †	32.7 (2.9) 52.5 (1.5) 58.8 (2.2)	38.8 (2.0) 58.4 (3.7) 57.2 (3.2)	14.7 (2.8) 28.8 (2.5) 	45.8 62.8 65.1* 63.3
	Ours	91.5 (1.8)	87.4 (0.5)	91.4 (0.6)	57.6 (5.5)	62.2 (6.0)	66.1 (4.1)	Ours	60.9 (1.9)	59.2 (3.0)	32.1 (2.0)	67.6
0	LOTClass [#] FT LM-BFF (man) LM-BFF (D-demo) KPT †	71.8 49.1 83.5 82.9 78.4	81.7 50.0 80.3 80.7 81.9	50.1 49.8 78.4 81.4 71.4	50.4 34.4 49.7 52.2 37.1	36.5 49.5 50.5 53.5 55.3	55.9 31.6 49.7 44.0 47.5	LOTClass [♣] FT KnPr KnPr (D-demo) KPT †	11.5 10.0 15.9 24.6	9.8 6.2 10.3 11.6	2.5 0.5 2.3 0.8	41.1 31.2 46.7 47.0* 45.7
	Ours	86.8	83.5	79.7	53.7	56.2	56.7	Ours	41.3	12.2	2.8	52.5

Analysis of Memorization

Definition of Memorization Measurement

we define memorization measures as to how the classification varies when a training instance z is deleted from the trainset. We define and derive the memorization score for a training instance z as follows :

Cross-domain Results

Model	Source	Target Domain		
	16-shot MR	SST-2	CR	
FT	76.9	71.4	64.7	
LM-BFF (man)	87.0	88.9	86.9	
LM-BFF (D-demo)	86.6	89.3	87.5	
KPT	86.8	86.8	86.7	
RETROPROMPT	88.0	91.4	88.8	
	16-shot QQP	MRPC	RTE	
FT	60.7	43.7	48.0	
LM-BFF (man)	65.4	20.9	65.5	
LM-BFF (D-demo)	68.2	38.8	66.2	
KPT	71.6	42.3	65.8	
RETROPROMPT	74.0	49.4	67.3	

Full-data Results





> Top-memorized Instances: Typical or Atypical?

we adopt SST-2 to analyze the memorization by judging the atypical of an instance by checking the percentage of positive phrases.

Mem Group		Negative		Postive			
	FT	LM-BFF	OURS	FT	LM-BFF	OURS	
Top-10%	34.29	32.78	30.23	68.75	69.71 86.39	75.67	
Bottom-10%	17.63	16.25	14.42	95.92	95.08	94.53	
		FT	LM-BFF		OURS		
MEM SCORE	MEM SCORE 4.597		0.1	21	0.032		