

# An Introduction to Prompting Methods

Yige Xu PhD Student, SCSE, NTU <u>yige002@e.ntu.edu.sg</u> <u>https://xuyige.github.io</u>

#### Outlines



- Development Line before Prompting
- Formulation of Prompting
- Prompt Learning and Training Strategies
- Verbalizer Learning
- Prompting Methods
- Recent Applications on Prompt

• Pre-trained Language Models



- Pre-train in large-scale unlabeled corpus and then fine-tune in "small-scale" downstream data have brought breakthrough on many NLP tasks
- It is hard to tune all parameters when PLM is becoming larger and larger
- Two solutions
  - Model compression and inference speedup, including knowledge distillation, quantization, parameter sharing, module replacing, and early exit
  - Parameter-efficient tuning methods, including Adapters, Part Parameter Tuning, and Prompt Tuning

- Adapters
  - Add the adapter module inside the Transformer layer and fix the PLM



- Adapters
  - Can be used in Cross-task/Cross-lingual/Cross-Modal transfer



• Part Parameter Tuning



- BitFit: Simple Parameter-efficient Fine-tuning for Transformer-based Masked Language-models (ACL 2022) Only the bias-term (or a subset of them) is fine-tuned
- Parameter-Efficient Transfer Learning with Diff Pruning (ACL 2021) Learn a diff vector and add this vector to the pre-trained model parameters
- Raise a Child in Large Language Model: Towards Effective and Generalizable Fine-tuning (EMNLP 2022) First find a subset of parameters and generate the gradient mask, then update part of the parameters based on the gradient mask

### **Formulation of Prompting**

- Pattern-Exploiting Training
  - M: Masked Language Model
  - L: A Set of Labels
  - V: Vocabulary
  - x: Input Sequence
  - [MASK]: Masked Token, in V
  - pattern: A function P that takes x as input and outputs a phrase or sentence P(x) in V\*
  - verbalizer: An injective function v: L->V
  - Pattern-verbalizer pair (PVP): (P, v)



#### **Formulation of Prompting**

#### • Notations



• Prompt Engineering



- Prompt engineering is the process to create a prompting function  $f_{\text{prompt}}(x)$  that maximize the performance on downstream tasks.
- In many previous works, one or more templates are searched for each task.
- Prompt Types:
  - Cloze Prompt: The language model is required to fill in the blanks of a textual string.
  - Prefix Prompt: The language model is required to generate the textual string by the given input as well as a prefix string.

- Manual Template
  - Manually define some patterns:
- p. Question: q? Answer: \_\_\_.
- p. Based on the previous passage, q? \_\_\_.
- Based on the following passage, q? \_\_\_\_. p
  - Advantages:
    - Help leverage the knowledge contained in pre-trained language model
    - Can be manual designed with prior expertise
    - Provide a perspective of better understanding the pre-trained language model

 $P_5(\mathbf{x}) = \dots$  News: a b

 $P_6(\mathbf{x}) = [$  Category: \_\_\_\_ ] a b

 $P_1(\mathbf{x}) = \_ a b \qquad P_2(\mathbf{x}) = a (\_ b b)$ 

 $P_3(\mathbf{x}) = \dots - a b \qquad P_4(\mathbf{x}) = a b (\dots)$ 

- Disadvantages:
  - Not flexible enough
  - Require much prior expertise



- $P_1(a) = \text{It was } a \quad P_2(a) = \text{Just } a$  $P_3(a) = a \text{. All in all, it was } a$ 
  - $P_4(a) = a \parallel$  In summary, the restaurant is \_\_\_\_\_

- Automated Template
  - Automatically search templates in discrete space (word index) or continuous space (word embedding)
  - Methods:
    - Discrete Prompts, include Prompt Mining (Jiang et. al. 2020), Prompt Paraphrasing (Yuan et. al. 2021, Haviv et. al. 2021), Gradient based Search (Shin et. al. 2020), etc.
    - Continuous Prompts. Include Prefix Tuning (Li and Liang, 2021), Tuning Initialized with Discrete Prompts (Zhong et. al. 2021, Shin et. al. 2020), Hard-Soft Prompt Hybrid Tuning (Liu et. al. 2021, Han et. al. 2021), etc.
  - Advantages:
    - Relax the constraint of embedding space
    - Remove the restriction that the template is parameterized by the PLM parameters
  - Disadvantages:
    - Can hardly be interpretable
    - Usually, cannot be easily transferred



- Training Settings
  - Zero-shot learning (PLM has been trained for LM to predict the probability, GPT-3)
  - Few-shot learning (Resource limitation scenarios, PET-TC)
  - Full-data learning (For best downstream task performance, PTR)
- Parameter Updating
  - No prompt, only fine-tuning: Only tuned LM params (BERT, RoBERTa)
  - Tuning-free prompting: Frozen LM params, no additional prompt params (GPT-3, BARTScore)
  - Fixed-LM prompt tuning: Frozen LM params, additional and tuned prompt params (Prefix-Tuning, WARP)
  - Fixed-prompt LM Tuning: Tuned LM params, with fixed additional prompt params (T5, PET-TC)
  - Prompt + LM tuning: Tune both LM params and additional prompt params (P-Tuning, PTR)



- Parameter Updating
  - No prompt, only fine-tuning: Only tuned LM params
    - Advantages: Simplicity, do not need prompt design. Can also fit to larger training datasets.
    - Disadvantages: LM usually overfit or unstable on small datasets.
    - Examples: BERT, RoBERTa
  - Tuning-free prompting: Frozen LM params, no additional prompt params
    - Also referred as "in-context learning": [TASK DESCRIPTION] x1 [SEP] y1 [SEP] x2 [SEP] y2 [SEP] x [SEP] [MASK] ?
    - Advantages: No parameter update process, which can be applied in zero-shot settings.
    - Disadvantages: (1) Requires much prompt engineering with prior expertise; (2) Slow at testing, cannot easily use large training datasets.
    - Examples: GPT-3, BARTScore



- Parameter Updating
  - Fixed-LM prompt tuning: Frozen LM params, additional and tuned prompt params
    - Advantages: Often outperforms tuning-free prompting, while retain knowledge in LMs and is suitable in few-shot scenarios (similar to tuning-free prompting).
    - Disadvantages: (1) Not applicable in zero-shot scenarios; (2) The learned representation has limitations in largedata scenarios (though effective in few-shot scenarios); (3) Prompts are usually hard to human-interpreted.
    - Examples: Prefix-Tuning, WARP
  - Fixed-prompt LM Tuning: Tuned LM params, with fixed additional prompt params
    - Advantages: Prompt or answer engineering is more specific to the task, allowing for more efficient learning.
    - Disadvantages: (1) Still requires manual prompt or answer engineering; (2) LM fine-tuned on one downstream task cannot be easily transferred to another downstream task.
    - Examples: T5, PET-TC
    - With null prompt, where directly concatenates "[X] [Z]" without any template words can sometimes achieves competitive accuracy.



- Parameter Updating
  - Prompt + LM tuning:
    - Tune both LM params and additional prompt params.
    - This settings is very similar to the standard pre-train and then fine-tune paradigm.
    - Advantages: Most expressive method, which is suitable for large dataset settings.
    - Disadvantages: Requires most computational resources. Have the similar limitations with pre-train + fine-tune paradigm (e.g., may overfit to small datasets)
    - Examples: P-Tuning, PTR



#### **Verbalizer Learning**

- Human-written Verbalizer
  - Manually define some verbalizers (Example: Yelp-5)

v(1) = terrible v(2) = bad v(3) = okayv(4) = good v(5) = great

- Automatic Verbalizer:
  - Use multiple verbalizers with a weight

$$q_{\mathbf{p}}(y \mid \mathbf{x}) = \frac{\exp\left(\frac{1}{|v_y|} \sum_{t \in v_y} M(t \mid P(\mathbf{x}))\right)}{\sum_{i=1}^k \exp\left(\frac{1}{|v_i|} \sum_{t \in v_i} M(t \mid P(\mathbf{x}))\right)}$$

- Automating Label Token Selection
  - Train a logistic classifier to predict the class label using the contextualized embedding of the [MASK] token
  - Substitute the hidden state with the MLM's output word embeddings to obtain a score
  - Construct the sets from the k-highest scoring words





#### **Verbalizer Learning**

- Knowledgeable Prompt-Tuning
  - Construct external knowledge base by Related Words, ConceptNet, WordNet, etc.
  - Remove the rare words





- Theoretical Support
  - Can Unconditional Language Models Recover Arbitrary Sentences? (NeurIPS 2019)
  - A PLM can generate a sentence through the identification of a point in the sentence space



Figure 1: We add an additional bias,  $W_z z$  (left, when  $\dim(z) \le d^*$ ) or  $Z = [z_1 \dots z_K]$  (right, when  $\dim(z) > d^*$ ), to the previous hidden and cell state at every time step. Only the z vector or Z matrix is trained during forward estimation: The main LSTM parameters are frozen and  $W_z$  is set randomly. In the right hand case, we use soft attention to allow the model to use different slices of Z each step.



- Prefix-Tuning
  - Prefix-Tuning: Optimizing Continuous Prompts for Generation (ACL 2021)
  - Add a prefix before the input and only tune the parameters from the prefix



Fine-tuning



- Prefix-Tuning
  - Prefix-Tuning: Optimizing Continuous Prompts for Generation (ACL 2021)
  - Add a prefix before the input and only tune the parameters from the prefix



#### Summarization Example

Article: Scientists at University College London discovered people tend to think that their hands are wider and their fingers are shorter than they truly are.They say the confusion may lie in the way the brain receives information from different parts of the body.Distorted perception may dominate in some people, leading to body image problems ... [ignoring 308 words] could be very motivating for people with eating disorders to know that there was a biological explanation for their experiences, rather than feeling it was their fault."

Summary: The brain naturally distorts body image – a finding which could explain eating disorders like anorexia, say experts.

#### Table-to-text Example

Table: name[Clowns] customerrating[1 out of 5] eatType[coffee shop] food[Chinese] area[riverside] near[Clare Hall]

Textual Description: Clowns is a coffee shop in the riverside area near Clare Hall that has a rating 1 out of 5 . They serve Chinese food .



- P-Tuning v2
  - P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks (ACL 2022)
  - Add prompts in each Transformer layer, which makes the model more stable





- LM-BFF
  - Making Pre-trained Language Models Better Few-shot Learners (ACL 2021)
  - Use T5 to generate the template automatically
  - Add some demonstrations for better comparisons



(c) Prompt-based fine-tuning with demonstrations (our approach)





- Pre-trained Prompt Tuning
  - PPT: Pre-trained Prompt Tuning for Few-shot Learning (ACL 2022)
  - Use three different tasks for prompt pre-training: sentence-pair classification, multiple-choice classification, and single-sentence classification
  - Can obtain a good initialization point of prompts



#### Pre-Training (Unlabeled Data) : Next Sentence Prediction

- LM-as-a-Service
  - Black-Box Tuning for Language-Model-as-a-Service (arXiv 2201.03514)
  - Wrap the input with the template and pass the wrapped input to the large-scale PLM
  - Use the result returned by the large-scale PLM as well as gradient-free methods to optimize the parameters from prompts





- LM-as-a-Service
  - Black-Box Tuning for Language-Model-as-a-Service (arXiv 2201.03514)





- Prompt for Data Augmentation
  - PromDA: Prompt-based Data Augmentation for Low-Resource NLU Tasks (ACL 2022)
  - Use prompt as well as the tag to help generating the input example, which can train a better prompt initialization point







- Prompt for Zero-shot Relation Extraction
  - RelationPrompt: Leveraging Prompts to Generate Synthetic Data for Zero-Shot Relation Triplet Extraction (Findings of ACL 2022)
  - Manually design the structured template for relation generator and relation extractor

Input	Relation: <label>.</label>
Example	Relation: Military Rank.
Output	Context: <sentence>. Head Entity: <subject>, Tail Entity: <object>.</object></subject></sentence>
Example	Context: Their grandson was Captain Nicolas Tindal. Head Entity: Nicolas Tindal, Tail Entity: Captain.
(a) Structured template for relation generator.	
Input	Context: <sentence>.</sentence>
Example	Context: Their grandson was Captain Nicolas Tindal.
Output	Head Entity: <subject>, Tail Entity: <object>, Relation: <label>.</label></object></subject>
Example	Head Entity: Nicolas Tindal, Tail Entity: Captain, Relation: Military Rank.

(b) Structured template for relation extractor.

- Vision-Language
  - A Good Prompt Is Worth Millions of Parameters: Low-resource Prompt-based Learning for Vision-Language Models (ACL 2022)







- Machine Translation
  - MSP: Multi-Stage Prompting for Making Pre-trained Language Models Better Translators (ACL 2022)



(a) Basic (single-stage) prompting for MT.



(b) Multi-stage prompting.



