# NeurIPS 2022 New Orleans. Nov 28-Dec 09, 2022



# Thirty-sixth Conference on Neural Information Processing Systems

# NL4Opt: Natural Language for Optimizing Modeling —— 3rd Place Solution

Jiayu Liu<sup>1,2</sup>, Longhu Qin<sup>1,2</sup>, Yuting Ning<sup>1,2</sup>, Tong Xiao<sup>1,2</sup>, Shangzi Xue<sup>1,2</sup>, Zhenya Huang<sup>1,2</sup>, Qi liu<sup>1,2</sup>, Enhong Chen<sup>1,2</sup>, Jinze Wu<sup>2,3</sup>

<sup>1</sup>Anhui Province Key Lab. of Big Data Analysis and Application, School of Computer Science and Technology & School of Data Science, University of Science and Technology of China, <sup>2</sup>State Key Laboratory of Cognitive Intelligence, <sup>3</sup>iFLYTEK AI Research (Central China), iFLYTEK Co., Ltd.

#### Introduction

#### **Background**

- Many real-world decision-making problems can be formulated and solved as mathematical optimization problems.
- ➤ Modeling a problem into a proper formulation is a complex and time-consuming process.

#### **Research Problem**

- > Can learning-based natural language interfaces be leveraged to convert linear programming (LP) word problems into a format that solvers can understand?
  - ➤ Challenge 1: detect problem entities from the problem description
  - ➤ Challenge 2: generate a precise meaning representation of the optimization formulation

### Subtask1 Summary

#### > Pre-processing

- > Data Augmentation
- > Swap variables, synonym replace for OBJ\_NAME, randomly replace numbers
- > Label Standardization
  - Label inconsistency: The word times after numbers are sometimes regarded as a part of PARAM and LIMIT entity, but sometimes not.
- entity

#### > Adversarial Training

Use Projected Gradient Descent[1] to improve the robustness and generalization ability of the model

- > xlm-roberta-base
- xlm-roberta-large

#### > Post-processing

## **Problem Definition**

### Subtask1

> LP problem description:

 $Q = \{w_0, w_1, ..., w_n\}$ 

#### Goal

Given

- $\triangleright$  LP problem entities:  $E = \{ E_0, E_1, ..., E_m \}$ •  $E_i = \langle I_i^{start}, I_i^{end}, c_i \rangle$ 
  - $I_i^{\text{start}} \in [0, n], I_i^{\text{end}} \in [0, n], c_i \in C$

#### Subtask2

#### Given

➤ LP problem description: Q > LP problem entities: E

#### Goal

- ➤ LP problem solution template
  - $A = \{S_0^{const}, S_1^{const}, ..., S_k^{const}, S^{obj}\}$

  - Declaration template :  $S = \{s_0, s_1, ..., s_n\}$

Remove times after numbers in PARAM and LIMIT

#### > Ensemble Learning

### Subtask2 Summary

#### >Pre-processing

> Repeat CONST\_DIR

#### >Prompt Re-design

- > Version1: Insert the prompt in the document
- Version2: Entity-enhanced prompt

#### >Encoder-decoder model

- > Model: BART
- > Adversarial Training
  - > A training method that introduces noise and improve the robustness and generalization ability of model

#### >Intermediate representation

**Ensemble** 

#### **Ensemble**

**>obj\_declaration**: version2 prompt-guided input **>const\_declarations**: version1 prompt-guided input + CONST\_DIR repeat + adversarial training

#### >Post-processing

- >Rule1: correct wrong operator direction
- >Rule2: correct wrong unseen numbers

object\_declaration

constriant declaration

constriant\_declaration

constriant\_declaration

# **Prompt Design**

#### >Pre-processing:

> Repeat CONST\_DIR

Avian would like the commercials to be seen by at least 86 million young girls and 72 million middle-aged women.

• 2 LIMIT entities are related to 1 CONST\_DIR. • 2 declarations are generated by 1 CONST\_DIR entity.

**Preprocessing** 

Avian would like the commercials to be seen by at least 86 million young girls and at least 72 million middle-aged women.

Repeat the CONST\_DIR entity before the second LIMIT entity

#### > Re-designed prompt

- ➤ Version1: Insert the prompt in the document
- ➤ Version2: Entity-enhanced prompt

<CONST\_DIR>at most</CONST\_DIR><s>A bakery sells two types of cakes. They sell a chocolate cake and a vanilla cake. Let's say they make x1 chocolate cakes, at a profit of \$10 each, and x2 vanilla cakes, at a profit of \$11 each (x1 and x2 are unknowns both greater than or equal to 0). The daily demand for these cakes is at most 50 chocolate cakes and at most 30 vanilla cakes. The bakery is short staffed and can make a maximum of 50 cakes of either type per day. How much of each cake should the bakery make in order to maximize profit?</s>

Original prompt of baseline

#### **Re-designed prompt**

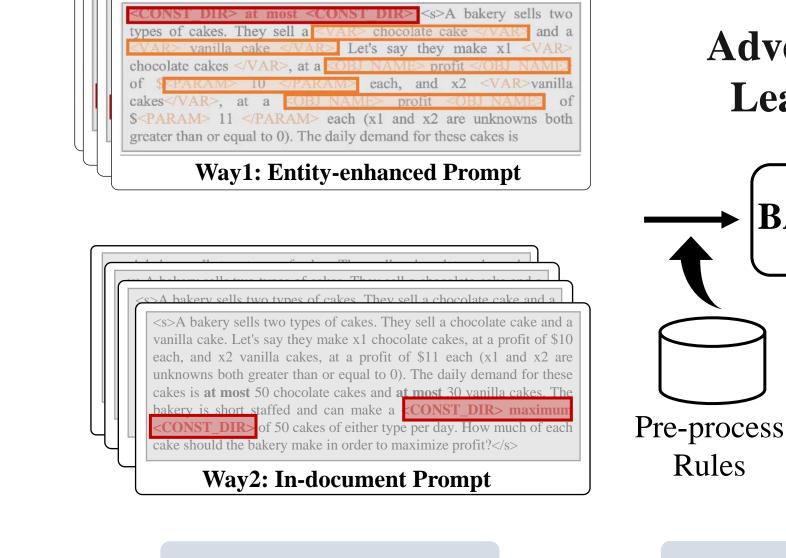
<s>A bakery sells two types of cakes. They sell a chocolate cake and a vanilla cake. Let's say they make x1 chocolate cakes, at a profit of \$10 each, and x2 vanilla cakes, at a profit of \$11 each (x1 and x2 are unknowns both greater than or equal to 0). The daily demand for these cakes is **CONST\_DIR**> at most **CONST\_DIR**> 50 chocolate cakes and at most 30 vanilla cakes. The bakery is short staffed and can make a maximum of 50 cakes of either type per day. How much of each cake should the bakery make in order to maximize profit?</s>

Version1: Insert the prompt in the document

<CONST\_DIR> at most <CONST\_DIR> <s>A bakery sells two types of cakes. They sell a <VAR> chocolate cake </VAR> and a <VAR> vanilla cake </VAR>. Let's say they make x1 <VAR> chocolate cakes </VAR>, at a <OBJ\_NAME> profit </OBJ\_NAME> of \$<PARAM> 10 </PARAM> each, and x2 <VAR>vanilla cakes</VAR>, at a <OBJ\_NAME> profit <OBJ\_NAME> of \$<PARAM> 11 </PARAM> each (x1 and x2 are unknowns both greater than or equal to 0). The daily demand for these cakes is <CONST\_DIR> at most <CONST\_DIR> <LIMIT> 50 </LIMIT>...

Version2: Entity-enhanced prompt

# Framework



**Prompt-guided Input** 

Rules

**Encoder-Decoder** 

**Adversarial** 

Learning

**Intermediate Representation** 

**Way1 Output** 

**Way2 Output** 

piect declaration

constriant\_declaration

object declaration

constriant declaration

Figure 1: Pipeline for our subtask 2 solution

Post-process

Rules

# **Datasets**

Dataset	#CONST_DIR	#LIMIT	#CONST_DECLARATION	
Train	2.27	2.53	2.78	
Dev	2.57	2.62	2.94	
Dataset	#OBJ_DIR	#OBJ_NAME	#OBJ_DECLARATION	
Train	1.0	2.96	1.0	
Dev	1.0	2.86	1.0	

Table1: Analysis about average number of declarations of each solution, along with related entities.

Processing Method	Train	Dev	
Add Duplicate Const	210	22	
Rule1	N/A	3	
Rule2	N/A	5	

Table2: Statistics about processed number of document by different pre-processing and post-processing for our subtask2 solution

# **Experiments**

Method	VAL_ACC	CONST_ACC	OBJ ACC	TEST_ACC
	_		<u></u>	_
baseline	0.6149	0.5502	0.7778	0.641
Version1 Prompt	0.7586	0.7149	0.8687	0.803
+ attack-fgm	0.7845			0.825
+ repeat CONST_DIR*	0.8649			0.856
Version2 Prompt*	0.6667	0.5542	0.9394	0.708
+ attack-fgm	0.6810			0.693
Ensemble				0.867
+ post-processing				0.880

Parser

a constraint with "GREATER OR FOUAL" operator

**Canonical** 

**Formulation** 

Table2: Main experiments results of our subtask2 solution

- Prompt-in method indicates the position of the declaration, significantly improving the performance.
- Span-in method brings in much entity information, greatly improving object declaration, but const declaration improve little, which may suffer from extra information burden.
- Repeating CONST\_DIR which is followed by two LIMIT alleviates the mismatch between the number of the ground truth and predicted declarations.

# Conclusion

- > We designed two powerful types of prompts. Experiments show that prompt design is important for text generation tasks, including what to prompt and where to prompt.
- We designed useful Post-processing and pre-processing methods based on mathematical knowledge, which directly correct some errors for model.
- We leveraged adversarial training to improve the robustness of the model, which helps for both Natural Language Understanding and Natural Language Generation.

### **Future work**

- > Data augmentation or more pre-processing rules.
- ➤ More attempts of prompt design.
- > Constraint the prefix of intermediate representation.



# **BDAA USTC**

Anhui Province Key Laboratory Big Data Analysis and Application