



## Problem Description

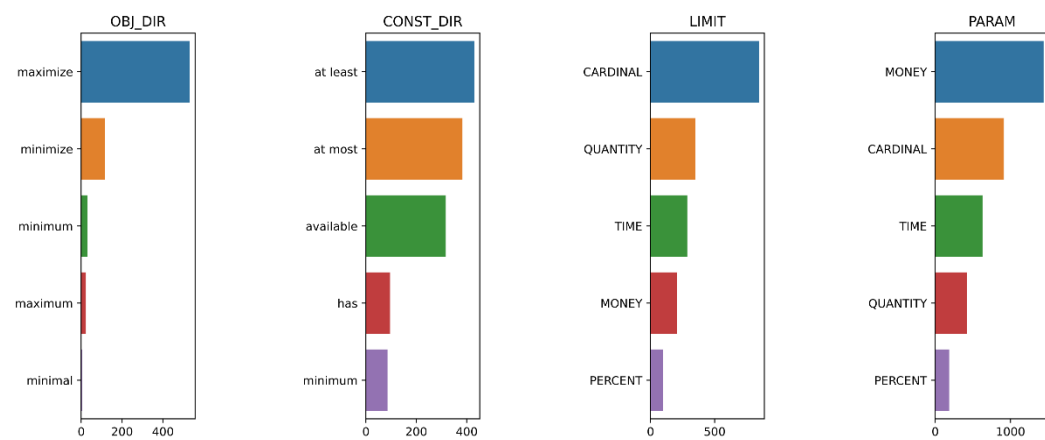
Given an expert formulated optimization problem in natural language, extract six named entities: **CONST\_DIR** (constraint direction), **LIMIT** (limit), **OBJ\_DIR** (objective direction), **OBJ\_NAME** (objective name), **PARAM** (parameter), **VAR** (variable). See example below:

The Notorious Desk company wants to promote a new brand of wine and wants to market it using a total market budget **CONST\_DIR** of \$ 87,000 **LIMIT**. To do so, the company needs to decide how much to allocate on each of its two advertising channels : ( 1 ) morning TV show **VAR** and ( 2 ) social media **VAR**. Each day, it costs the company \$ 1,000 **PARAM** and \$ 2000 **PARAM** to run advertisement spots on morning TV **VAR** show and social media **VAR** respectively. The expected daily reach **OBJ\_NAME**, based on past ratings, is 15,000 **PARAM** viewers for each morning show **VAR** spot and 30,000 **PARAM** internet users for a social media spot **VAR**. The chief marketer knows from her experience that both channels are key to the success of the product launch. She wants to plan at least **CONST\_DIR** 4 **LIMIT** but no more than **CONST\_DIR** 7 **LIMIT** morning show **VAR** spots. In addition, the social media spots **VAR** needs to be at least **CONST\_DIR** 30 **LIMIT** due to pricing tier policy. How many times should each of the media channels be used to maximize **OBJ\_DIR** the reach **OBJ\_NAME** of the campaign ?

## Data Characteristics

Number of Samples: Training - 713, Dev - 99

Are there any frequently occurring key-phrases or themes in these entities?



Common key-phrases in **CONST\_DIR** and **OBJ\_DIR** and Common Themes in **LIMIT** and **PARAM**

Are there any patterns in Structure and Verbiage in these entities?

**Conjuncting Noun Chunks**

A factory in India produces rice **VAR** and corn **VAR**.

Firefighting units can either send units of firefighters **VAR** or volunteer fire patrols **VAR**.

**Conjuncting Prepositional Chunks**

There are three types of commercials. Commercials with famous actors **VAR**.

commercials with regular people **VAR**, and commercials with no people **VAR**.

**Hyphens**

A clothing company makes blue **VAR** and dark blue t-shirts **VAR**.

**Quotes**

An MCA checks a patient's eye pressure one - by - one either by using a tonometer **VAR** or a "puff of air" test **VAR**.

Common Patterns in **VAR**

**Pattern Extraction**

```
# use all keywords tagged as OBJ_DIR in training data
OBJ_DIR_KEYWORDS = ["maximize", "minimize", ...]
OBJ_NAME_PATTERNS = [
    (<LNAME>: (<1st>: OBJ_DIR_KEYWORDS)),
    (<POS>: (<1st>: (<DET>, "PART", "ADJP"), <OP>: ">"),
    (<POS>: (<1st>: (<PART>, "ADJ"), <OP>: ">"),
    (<POS>: "NOUN", <OP>: ">"),
    (<POS>: "VERB", <OP>: ">"),
    (<POS>: "NOUN", <OP>: ">"),
]
```

**Subject Object Extraction**

I want my profit **SUBJECT** to be maximized **OBJ\_DIR**.

I want the one with minimal **OBJ\_DIR** cost **OBJECT**.

Common Patterns in **OBJ\_NAME**

## Experimental Protocol

### Feature Engineering

- ❑ CRF model exploring basic grammatical and morphological features
- ❑ CRF model exploring grammatical, morphological and engineered features inspired from the Data Characteristics

### Feature Learning

- ❑ Token-classification model using RoBERTa large
- ❑ Ensemble of two separate token-classification models one for just OBJ\_NAME and VAR and the other for the rest
- ❑ Token-classification model with a modified cost function to optimize for mistakes in OBJ\_NAME and VAR
- ❑ Token-classification model using XLM-RoBERTa and curriculum learning
- ❑ Token-classification model using XLM-RoBERTa fine-tuned on Optimization Corpora

### Hybrid

- ❑ CRF Model combining best performing Feature Engineering and Feature Learning techniques

## Augmentation Strategies

### Up sampling via Duplication of in-frequent patterns

- ❑ **OBJ\_DIR** is generally a verb (e.g., maximize, minimize) but there are a few examples, where **OBJ\_DIR** is also an adjective (e.g., I want the cost to be minimal)
- ❑ **VAR** is mostly a Conjuncting noun chunk. Conjuncting prepositional phrases are an infrequent pattern (e.g., He does commercials with famous actors and commercials with regular actors)
- ❑ **OBJ\_NAME** is **OBJ\_DIR** followed by a noun phrase / prepositional phrase. **OBJ\_DIR** followed by multiple prepositional phrases is a rare pattern (e.g., maximize the number of action figures; minimize the number of batches of cookies)

**Augmenting Last Two Sentences:** In most cases, for **OBJ\_NAME** tokens to be tagged correctly it is imperative that the objective is known first. For example:

A doctor can prescribe two types of medication for high glucose levels, a diabetic pill **VAR** and a diabetic shot **VAR**. Per dose, diabetic pill **VAR** delivers 1 **PARAM** unit of glucose reducing medicine and 2 **PARAM** units of blood pressure reducing medicine **OBJ\_NAME**. Per dose, a diabetic shot **VAR** delivers 2 **PARAM** units of glucose reducing medicine and 3 **PARAM** units of blood pressure reducing medicine **OBJ\_NAME**. In addition, diabetic pills **VAR** provide 0.4 **PARAM** units of stress and the diabetic shot **VAR** provides 0.9 **PARAM** units of stress. At most **CONST\_DIR** 20 **LIMIT** units of stress can be applied over a week and the doctor must deliver at least **CONST\_DIR** 30 **LIMIT** units of glucose reducing medicine. How many doses of each should be delivered to maximize **OBJ\_DIR** the amount of blood pressure reducing medicine **OBJ\_NAME** delivered to the patient ?

**Pseudo Label Data generation:** Use paraphrase corpora like WordNet and PPDB to generate pseudo label data

## Hybrid Model

### Feature Engineering

- ❑ Grammatical Features
- ❑ Morphological Features
- ❑ Gazetteer Features
- ❑ Features exploiting syntax and verbiage
- ❑ Features are extracted at each word position and a window around it

### Feature Learning

- ❑ Label predictions from a trained RoBERTa transformers model
- ❑ Large variant and the base variants is used for comparison

### Conditional Random Field

## Selected Results

Model Name	CONST_DIR		LIMIT		OBJ_DIR		OBJ_NAME		PARAM		VAR		Average Micro F1 (Dev)	Average Micro F1 (Test)
	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall	Precision	Recall		
Grammatical and Morphological Features + CRF	0.956	0.854	0.904	0.954	0.979	0.929	0.649	0.353	0.958	0.916	0.795	0.714	0.816	-
Grammatical, Morphological, Gazetteer, Structural Features + CRF	0.960	0.858	0.931	0.942	0.990	0.970	0.726	0.544	0.953	0.935	0.823	0.787	0.853	-
RoBERTa Large	0.895	0.902	0.984	0.950	0.990	1.000	0.668	0.597	0.965	0.983	0.916	0.940	0.904	-
[Infrequent Pattern Upsampling] + RoBERTa Large	0.947	0.909	0.984	0.950	0.990	0.990	0.628	0.615	0.961	0.979	0.906	0.947	0.903	-
Pre-trained XLM-RoBERTa Large on Textbooks	0.901	0.897	0.987	0.953	0.989	0.999	0.665	0.583	0.971	0.989	0.918	0.946	0.907	
[Last Two Sentence Augmentation] + Grammatical, Morphological, Gazetteer, Structural Features + RoBERTa Model Predictions + CRF + RandomSearchCV	0.946	0.890	0.980	0.942	0.990	1.000	0.730	0.668	0.957	0.983	0.935	0.953	0.919	0.920

## Discussions & Observations

Scope for Aleatoric Uncertainty - Similar sequences annotated differently in Train and Dev

How should the bakery operate to maximize <b>OBJ_DIR</b> total profit <b>OBJ_NAME</b> ?
How many of each type of transportation should the company schedule to move their lumber to minimize <b>OBJ_DIR</b> the total cost <b>OBJ_NAME</b> ?
How many of each type of donut should be bought in order to maximize <b>OBJ_DIR</b> the total monthly profit <b>OBJ_NAME</b> ?
If the chemical company needs to make at least <b>CONST_DIR</b> 900 <b>LIMIT</b> au of the acidic liquid and 1200 <b>LIMIT</b> au of the basic liquid per minutes <b>OBJ_NAME</b> should each reaction be run for to minimize <b>OBJ_DIR</b> the total time <b>OBJ_NAME</b> needed ?
How many of each should the pharmaceutical manufacturing plant make to minimize <b>OBJ_DIR</b> the total number of minutes needed <b>OBJ_NAME</b> ?
Cautious Asset Investment has a total <b>CONST_DIR</b> of \$ 150,000 <b>LIMIT</b> to manage and decides to invest it in money market fund <b>VAR</b> , which yields a 2 % <b>PARAM</b> return <b>OBJ_NAME</b> as well as in foreign bonds <b>VAR</b> , which gives and average rate of return <b>OBJ_NAME</b> of 10.2 <b>PARAM</b> %.
To do so, the company needs to decide how much to allocate on each of its two advertising channels : ( 1 ) morning TV show <b>VAR</b> and ( 2 ) social media <b>VAR</b> . Each day, it costs the company \$ 1,000 <b>PARAM</b> and \$ 2000 <b>PARAM</b> to run advertisement spots on morning TV <b>VAR</b> show and social media <b>VAR</b> respectively.

Excerpts from Train (green) and Dev (yellow) highlighting annotation inconsistency for similar sequences