

Natural Language Processing HW2

Anonymous ACL submission

Abstract

This document provides details about the second NLP homework assignment. The tasks involve initially fine-tuning and evaluating a large language model (LLM). Following this, I will augment the training data to enhance the model's robustness and improve its performance.

1 Model and dataset

Regarding the model used, I decided to load the pre-trained distilbert-base-uncased model and subsequently fine-tune it with the Fever dataset provided in the homework slides.

1.1 Model and parameters

As mentioned, my homework loads the pre-trained distilbert-base-uncased model, a faster and smaller version of the BERT model, to allow for training and evaluation in reasonable timeframes while still achieving decent results. The parameters used for training are: batch_size = 32, learning_rate = 1e-4, weight_decay = 0.001, and epochs = 1.

1.2 Dataset

Speaking of the dataset, I worked with a subset of FEVER, loaded from HuggingFace using the datasets library. This subset contains samples for training, validation, and testing, with the following features: 'id', 'premise', 'hypothesis', 'label', 'wsd', and 'srl'. The steps I performed to work with this dataset are as follows:

- **Load the dataset;**
- **Map the labels of all subsets (0 = Entailment, 1 = Neutral, 2 = Contradiction);**
- **Tokenize the dataset:** using the tokenizer provided by distilbert-base-uncased, it assigns a numerical value (token) to each word in the premise and hypothesis, enabling the model to process the data.

2 Data augmentation

Regarding the data augmentation part, utilizing the "wsd" (Word Sense Disambiguation) and "srl" (Semantic Role Labeling) features, along with the assistance of pre-existing NLP models for data processing, I developed and implemented five methodologies, described below.

2.1 Hypernym/Hyponym replacement

The first method I devised involves replacing a word (specially a Noun) in the sentence with its hypernym or hyponym based on a prediction made by the pre-trained BERT model. This is done by masking the target word and predicting the best replacement word within the sentence. Specifically, the steps followed are:

1. Create a list of candidates, obtained from the method that returns the hypernyms and hyponyms of the lemma of the target word.
2. Mask the target word in the sentence, tokenize the sentence with the mask, and process it through the model.
3. Use the model to generate a list of 100 possible tokens to complete the sentence.
4. Calculate the cosine similarity between the generated tokens and all the candidates in the list from step 1.
5. Finally, select the word with the highest similarity to the hypernyms/hyponyms and insert it into the original sentence, converting the word to plural using the inflect library, if necessary.

2.2 Agent and patient swapping

The second technique involves identifying words with the roles of Agent and Patient in the sentences and simply swapping them to generate a sentence that is often false relative to the initial hypothesis.

072 Nonetheless, a method was created to calculate the
073 similarity between the generated sentence and the
074 original one. Based on a certain threshold, the label
075 of the sentence is appropriately changed.

076 **2.3 Antonym replacement**

077 The third methodology involves swapping a spe-
078 cific word (adverb or adjective) with its antonym
079 (obtained using the method that returns antonyms
080 through WordNet of a lemma), and accordingly
081 changing the assigned label.

082 **2.4 Paraphrasing**

083 The fourth method exclusively utilizes a pre-trained
084 model, specifically Pegasus, tasked with taking
085 the hypothesis from the dataset and paraphras-
086 ing it while retaining the meaning. Specifically,
087 the model is asked to return three paraphrased
088 sentences, and through a similarity comparison
089 method, the paraphrase with the highest similarity
090 score to the original sentence is chosen to preserve
091 the same meaning.

092 **2.5 Proper Noun replacement**

093 The last method I conceived involves replacing
094 proper nouns with other proper nouns that are as
095 similar as possible. This technique follows several
096 steps:

- 097 1. Gathering lists of proper names from external
098 sources such as GitHub or specific libraries,
099 including names of people (male or female)
100 and cities.
- 101 2. Creating a replacement dictionary where the
102 key is the original proper noun in the sentence,
103 and the value is the proper noun with the high-
104 est similarity score, obtained by comparing
105 the initial name with the lists gathered in step
106 1.
- 107 3. Implementing a method to substitute proper
108 nouns in the hypotheses of the sentences, with
109 appropriate changes to the label as a result.

110 **3 Reasons of the methods and Data 111 generation**

112 Among the listed methods, I focused more and
113 found more challenging to implement the one re-
114 lated to replacing with hypernyms and hyponyms,
115 as well as the one concerning proper nouns. Specif-
116 ically for the first method, simply substituting

117 through hypernyms and hyponyms led to person-
118 ally unsatisfactory results. For this reason, I tried
119 to adapt a strategy that allowed me to draw on these
120 two concepts but ultimately resulted in replacing
121 the word with something rather similar, rather than
122 the exact lemma derived from WordNet. This ap-
123 proach inevitably introduces ambiguity and errors,
124 especially regarding coherence and similarity be-
125 tween the original and the new sentence, as many
126 factors come into play. Nevertheless, my task was
127 also to try to minimize such errors by generating
128 sensible sentences.

129 Finally, for the automatic and random generation
130 of the new dataset, I initially randomly selected
131 15,000 samples from the initial training set. Subse-
132 quently, I developed a method that randomly chose
133 the function and methodology to apply to each pre-
134 viously selected sample. This approach allowed
135 me to create a new dataset consisting of 15,000
136 samples generated randomly (plus the other 50,000
137 from the original set) through augmentation, ready
138 to be used for fine-tuning the model.

139 **4 Results and comparison**

140 Regarding the results obtained, the first fine-
141 tuned model using the original datasets returned:
142 'eval_accuracy': 0.7022, 'eval_f1': 0.6860. These
143 values, while not excellent, reflect modest ar-
144 chitectural choices given the limited computa-
145 tional resources at my disposal. With more pre-
146 cise decisions, such as selecting a better training
147 epoch or a more powerful and efficient model like
148 RoBERTa, significantly better results could have
149 been achieved.

150 As for the adversarial set, the results
151 are also suboptimal, reported as follows:
152 'eval_accuracy': 0.5014836795252225, 'eval_f1':
153 0.5038331002848744.

154 After performing data augmentation, fine-tuning
155 was painstakingly carried out on the same model
156 using the new dataset. Considering the architec-
157 ture and parameter choices, the results showed
158 only slight improvement (at least they did not
159 worsen), and are likely much more improvable
160 with more accurate decisions. Specifically, I
161 achieved: 'eval_accuracy': 0.7092260603410582,
162 'eval_f1': 0.692296233672235. Regarding the
163 adversarial set provided in the slides, I obtained:
164 'eval_accuracy': 0.5133531157270029, 'eval_f1':
165 0.5140298545393321.

```
DatasetDict({
    train: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label', 'wsd', 'srl', 'input_ids', 'attention_mask'],
        num_rows: 51086
    })
    validation: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label', 'wsd', 'srl', 'input_ids', 'attention_mask'],
        num_rows: 2288
    })
    test: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label', 'wsd', 'srl', 'input_ids', 'attention_mask'],
        num_rows: 2287
    })
})
```

Figure 1: Original dataset structure.

```
My augmented dataset:
DatasetDict({
    train: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label'],
        num_rows: 61086
    })
    validation: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label'],
        num_rows: 7288
    })
    test: Dataset({
        features: ['id', 'premise', 'hypothesis', 'label'],
        num_rows: 2287
    })
})
```

Figure 2: Augmented dataset structure.