

Report

ML Problem Description and Potential Approaches :

Given a dataset of inputs, i.e. texts/tweets, and the emotion label conveyed by each tweet among 4 labels (anger, optimism, sadness, joy), create a classifier that can map the inputs to their outputs well.

Generally speaking, there are two approaches to the this problem:

1. **Finetuning LLM-based models (FT)** to achieve the desired result: Given that we have 1000+ data points and 4 labels, this problem could be solved practically by leveraging the generalization abilities LLMs have as they have been trained on large amount of text and they perform language modelling pretty well.

To perform transfer learning, we can simply optimize the full set of model parameters while replacing the classification head with a new randomly initialized one that specializes in our task of interest, i.e. predict one of the 4 categories.

2. **Prompt-based zero/few shot in-context learning (ICL) with text generative LLMs:** text-generation LLMs have shown to perform incredible zero-shot learning capabilities across many text-based tasks and therefore, we could experiment with performing in-context learning. In particular, we could instruct the model via prompting to generate the correct class of the question-answer pair without the requirement for any finetuning or model parameter optimization.

Research findings about 1 vs. 2:

- *Data hunger:* While ICL does require minimal to zero amount of data, FT needs a high-quality good amount of data for the optimization to generalize well to unseen examples. [1,2,3]

- *Prompt engineering*: While FT does not rely heavily on prompt formatting, ICL is shown to be quite sensitive to instructions and prompts that are essentially responsible for instructing the LLM to perform the task. In addition, it has been shown that the accuracy of prompt-based few-shot classification can vary to a large degree depending on the examples being provided in the prompts. [1,2]

- *Model family*:

- FT makes use of encoder-only transformer architectures like BERT which is known for achieving SOTA performance on language understanding tasks like text classification as they have been trained to optimize for masked language modelling (MLM), leading to rich language representations that can be used across different downstream tasks. [1,2,3]

- In addition, text generative LLMs, which have decoder-only architectures, like LLAMA and GPT families, can be finetuned in the same manner as Bert-like models by relying on the features extracted from the last hidden layer after decoding the last token in a sequence, followed by a new classification head. [1,2]

- For zero/few shot learning, decoder-only LLMs are used as they have been trained for text generation and they excel at following instructions and predicting next relevant words. [1,2,3]

- *Performance*:

- FT with encoder-only models (small models relative to decoder-only LLMs) still outperform all other methods in text classification across different tasks requiring good amount of data.

As a rule of thumb, FT works well on 500+ examples and the more data the better. [1]

- While applying FT to decoder-only LLMs is shown to achieve comparable good results to Bert-like models' performance, training decoder-only is usually very compute-expensive as they have 10x or 100x or 1000x parameters than Bert-like models. In that case, we can only use PEFT

methods to optimize their parameters, which can be suboptimal or achieve results similar to encoder-only models at best. [1,2]

- For zero/few shot ICL, only super large LMs can achieve comparable results to SOTA Bert-like models [1,2,3]. This makes sense as the larger the model the more powerful and generalizing it is.

It has been shown that using the instruct-finetuned version of these models is way better than using the pretrained versions in both FT and ICL. [1,2]

Experiments:

Note: All FT experiments are shown on WandB [here](#). Some specific information about each experiment is attached to the description of each experiment in the overview section!

Evaluation:

Since it is a classification task and we have an imbalanced dataset, I relied on calculating the F1-score with a 'weighted' mode to account for the importance of the class given its real representation in the dataset (regardless of train, val or test as they have the same class distribution through stratified splitting). This is motivated by the fact that the dataset class distribution represents real-world production case.

In addition to the F1-score, I log weighted recall, precision, the validation cross entropy loss and the per-label F1 score to WandB for further investigation.

FT with Bert-like models:

Given all the findings from literature, I went for experimenting with Bert-like models as we have good amount of data to make use of and they are relatively small models so I can finetune their full set of parameters.

I experimented with several hyperparameters to get best performance, namely learning rate, batch size, gradient clipping, weight decay, and dropout. In addition, I experimented with several pretrained architectures such as BERT, BERTweet, and more.

I found the following hyperparameters as best configuration in [this experiment](#) achieving *weighted_F1=83.5% on the validation set*:

```
model_name="vinai/bertweet-base"  
bsize=64  
lr=0.00003  
weight_decay=0.01  
max_grad_norm=1
```

FT with decoder-only models:

In the notebook, I implemented the code to be able to experiment with larger models to see if finetuning them would improve the previous results achieved by Bert. In particular, LLAMA 3-8B, which is known to be trained on a super large dataset, can be finetuned using LoRA for finetuning.

Zero/few shot ICL:

From the research findings above, I found that using super large models can benefit ICL pretty much so I chose to go for one of the largest LLMs open-sourced which is instruct-LLAMA 3-70B. I used 4-bit quantization to be able to use it on an A100-80GB gpu for inference.

From [1], we see that prompt design is really important for achieving good performance so I decided to implement three kinds of prompts to check which one works the best. I experimented with few shot ICL after I optimized the prompt design in zero-shot setting.

1. Generic prompt:

```
Return the {target} {labels} with one word/label response without any additional text.  
Answer: {target_with_brackets}.
```

where target is the target word to predict, i.e. in our case it is the type of emotion, labels is the list of labels among which the model shall choose, and the target_with_brackets is just {target} so it understands that it should be one-worded answer, which enforces its formatting to be as brief as possible.

2. Generic prompt + task-specific info which is:

```
""""This tweet is posted on Twitter and the text expresses
one of the following emotions:
- `anger`: Feelings of frustration or hostility, often exp
ressed through irritation or outrage.
- `joy`: Positive feelings of happiness or pleasure, shown
through excitement or delight.
- `optimism`: Hopeful and positive outlook, indicating enc
ouragement or belief in good outcomes.
- `sadness`: Feelings of sorrow or disappointment, express
ed through a sense of loss or melancholy.
""""
```

3. 1+2+ more context for the model:

```
"This tweet is posted on the social media platform Twitter
and you should classify it based on the emotion this tweet
expresses."
```

When performing zero shot ICL with prompt type. 1, I get on the validation set:

```
{'f1': 0.6010573614719079,
'precision': 0.627894781090178,
'recall': 0.5828877005347594,
'anger': 0.657439446366782,
'joy': 0.6428571428571429,
'optimism': 0.12345679012345678,
'sadness': 0.6043956043956044}
```

When performing zero shot ICL with prompt type. 2, I get on the validation set:

```
{'f1': 0.805891016161514,
'precision': 0.8139610480167755,
'recall': 0.7994652406417112,
```

```
'anger': 0.8974358974358975,  
'joy': 0.84375,  
'optimism': 0.36923076923076925,  
'sadness': 0.7374301675977654}
```

When performing zero shot ICL with prompt type. 3, I get on the validation set:

```
{'f1': 0.8019738480059521,  
'precision': 0.8082388143202798,  
'recall': 0.7967914438502673,  
'anger': 0.89171974522293,  
'joy': 0.8272251308900523,  
'optimism': 0.3492063492063492,  
'sadness': 0.7555555555555555}
```

Closely aligning with literature's findings, adding explicit task info and context about the task to the instruction prompt boosts the performance significantly, which is comparable to the SOTA bert-like models' performance.

Choosing prompt type 2. as our default prompt, I proceeded to the few shot ICL experimentation. By randomly selecting examples to be included in the prompt, I run 1-shot ICL 3 times and got various F1 scores on the validation set:

Run 1:

```
{'f1': 0.7839623191027408,  
'precision': 0.7944503473895672,  
'recall': 0.7780748663101604,  
'anger': 0.8372093023255814,  
'joy': 0.8181818181818182,  
'optimism': 0.4,  
'sadness': 0.7717391304347826}
```

Run 2:

```
{'f1': 0.8036112995286697,  
'precision': 0.812680300621477,
```

```
'recall': 0.7967914438502673,  
'anger': 0.8896103896103896,  
'joy': 0.8426395939086294,  
'optimism': 0.3384615384615385,  
'sadness': 0.7528089887640449}
```

Run 3:

```
{'f1': 0.7934341588685907,  
 'precision': 0.796626043059714,  
 'recall': 0.7914438502673797,  
 'anger': 0.8451612903225807,  
 'joy': 0.8383838383838383,  
 'optimism': 0.4406779661016949,  
 'sadness': 0.7624309392265194}
```

As we see, few-shot ICL does not improve zero-shot ICL on average and still achieves SOTA F1 scores compared to our finetuned Bert-like baseline.

By running 2-shot ICL, I did not see drastic improvement of the F1 score compared to one-shot ICL.

From [1], it seems that the choice of examples in the prompt influences the performance drastically. Therefore, to generalize better and have a better strategy for choosing examples, I have implemented an idea which is as follows:

- Given a new input, i.e. a tweet, we can retrieve from the training set the most similar tweet to that input across all labels based on Bert-embedding similarities. Then, we can introduce these tweets with their labels, which are most matching to our input tweet, in the instruction prompt and perform few-shot ICL.

Another heuristics to try out for example selection is [here](#).

Takeaway

Taking all results together, we can see that these findings are very reliable as they align with research findings about LLM-based text classification. The solution to

this problem has different approaches to follow but since we have 1000+ data points, we should rely on the full finetuning of relatively small Bert-like models as we can get more robust task-specific SOTA classifiers.

For production use cases, we should solve the classification problem in general as follows:

1. In case we have 1000+ data points, we should rely on the full finetuning of relatively small Bert-like models as we can get more robust task-specific SOTA classifiers.
2. In case we have less data, the results of FT can become less reliable and therefore, we should rely on few-shot ICL. The experiments I did show that few-shot ICL works very well and achieves SOTA results if prompt design and the choice of LLM is done properly.

In particular, we should rely on super large models like GPT-4o or LLAMA-3-70B and use a prompt design that includes the task information, context about the classification problem and the same format as proposed above. Also, we should follow the heuristics I implemented for choosing the examples to be introduced in the prompt, which makes the text generation process more relevant to the input to be classified and more stable for achieving high accuracy.

