

NLP Summarization Using Multiple Languages

By: Dylan Palmer

Abstract

As an attempt at making Amazon reviews more readable for consumers, I followed huggingface's Summarization project and tried to increase the accuracy of their results. While there was some success in changing hyperparameters, most trials ended in failure due to hardware limitations or given parameters being already optimized.

1 Introduction

Summarization is a task used for taking long or complicated documents and finding the most useful or the most relevant information from said documents, then condensing that into a short, readable form. This is especially used to help make dense, domain specific knowledge like research papers or medical papers more accessible. Summarization is also useful for making any other knowledge more digestible.

I am trying to make it easier to extract relevant information from product reviews. It is difficult to sift through products and reviews at the best of times; by simplifying reviews, it allows someone to get a general consensus on a product. To achieve this result, I am using the huggingface Summarization task code. [5]

2 Related Work

There are multiple types of summarization: extractive, and abstractive. Extractive summarization finds the most important words of the original document or text then strings them together as the summary. Abstractive summarization creates new phrases based on the original text using trained models to predict relations between original texts and their

summaries. This is more difficult than extractive summarization and therefore runs into more problems because of that.

With the rise in transformers, the rise in abstractive summarization has followed.[1] These researchers experimented with ways to increase the effectiveness of summarization. They found attention to be better than bag-of-words, though one researcher found that combining the old, extractive, bag-of-words method with the "new" attention increased accuracy.[2]

3 Model

3.1 Math

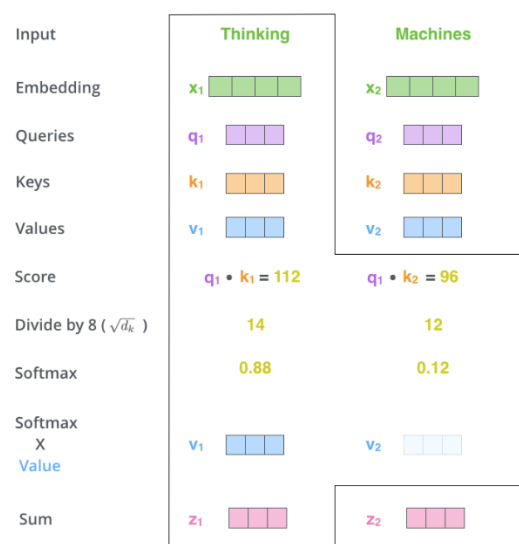


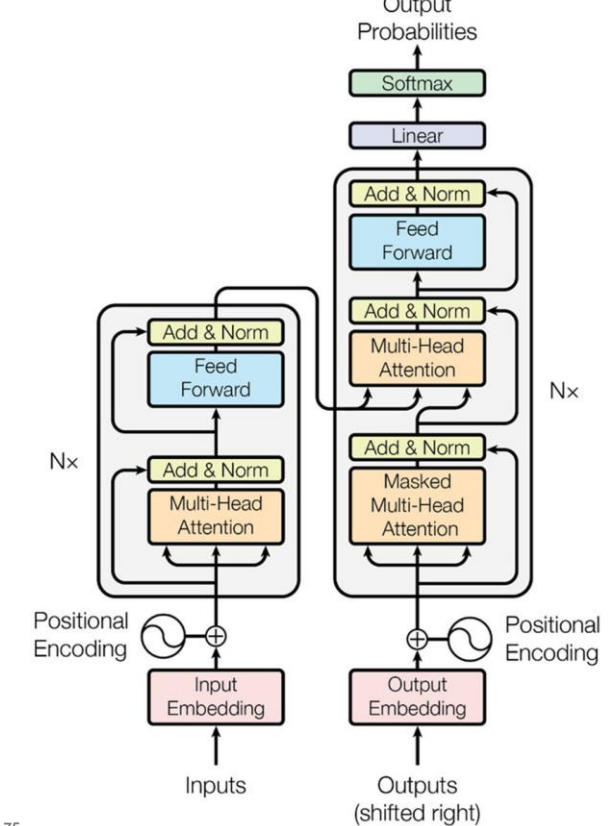
Figure 1: Attention equation [0]

by multiplying the initial input with the three matrices that are updated during training. These vectors are called query, key, and value vectors. Next, multiply the query and key, divide that number by the square root of d_k , then apply a softmax to get a most likely outcome for what is important. Multiply this by the value vector for the final result.

71

This attention picks out specific words or features that are important and passes that on to the feed forward network as shown below

72



75

Figure 2: Transformer Layout [3]

76

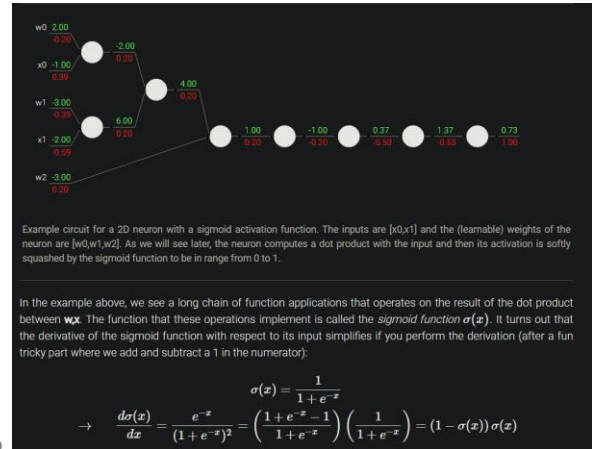
77

When using multi-head attention, each “head” performs the attention process, hopefully picking out different features from the encoding to pass on.

82

Backpropagation is an important operation for any machine learning network. It involves taking a function, one of several are commonly used, and feeding the result of the function to the last layer, and successive derivatives to each layer before that.

88



89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

Figure 3: backpropagation [6]

By having weights that update backwards it helps train the layers faster.

3.2 Model and Preprocessing

This model is using the multilingual Amazon review dataset, which consists of 200,000 training, 5,000 test and 5,000 validation samples in five languages. These reviews span all categories from Amazon products, from books to beauty. I will be focusing on using two languages, English and Spanish.

As it takes too long to process all 200k samples, I first stuck both the English and Spanish datasets into datagrams, then filtered them by category to get a subsample. I concatenate, then shuffle these two subsamples to get a single mix of English and Spanish. Lastly, I remove titles that are too short for better results.

After finishing setting up the data, next is the tokenizer. As I am using the mt5-small model, I also use the mt5-small tokenizer. I create a function to tokenize both the body and the title of the reviews and map it to the processed data. Next, I create the model using mt5-small preprocessed model. MT5 is used due the need to train multiple languages. The small version is used to save on GPU memory and make training times faster. The last things needed are setting up the necessary arguments such as learning rate and weight decay, using a data collator, and finally, using a sequence-to-sequence trainer.

124

125

126 **4 Experimental Results**

127 **4.1 Setup**

128 As all values had already been chosen for the initial
129 run, my first job was only getting the code to train,
130 then getting it to train on GCP. Only then did I start
131 trying to improve the results.

132 **4.2 Process**

133 Converting the code to run on GCP was a matter of
134 understanding and switching some jupyter
135 notebook conventions to work with a python script.
136 As there were installs inside the script I
137 simply pre-installed those packages in GCP and
138 removed them from the code. Since the code was
139 from huggingface, it used the huggingface site to
140 save and load data about the models trained. While
141 I could have gotten this to work on GCP, it was
142 easier to get the model to run without using the
143 huggingface hub at all.

144
145 Since all initial values were picked, my task
146 became changing each hyperparameter to see if
147 that gave better results. Trying mt5-base caused the
148 GPU to run out of memory. So did moving from a
149 batch of 8 to 16. I tried using gradient
150 accumulation to combat this error, but either I was
151 using it wrong, or it just didn't help. Both the
152 ROGUE and the loss were worse. Something else
153 I noticed is the sample rate was only a quarter of a
154 normal run. I assume this is due to it only updating
155 the gradient every four runs. To combat this, I
156 increased the epochs both with and without
157 gradient accumulation. The result was a worse
158 ROGUE and the loss again for both. I was later
159 informed that early stopping code is needed for
160 higher epoch numbers, but that's not something I
161 learned how to do. I also never figured out how to
162 zero out the weights of a pretrained model, so I do
163 not have a baseline.

165 **4.3 Results**

166
167 The only thing that did work was increasing the
168 sample size of the subsample for training. When
169 trying the full dataset, the training time was 32 hrs,
170 an unreasonable number, so I just added another
171 category to the filtering. This caused the original
172 ROGUE-L to go from 17 to 18.5

Run	ROGUE-L	Loss
Baseline	17.1	3.03
40 - Epoch	18.45	3.13
Extra Samples	18.56	2.63

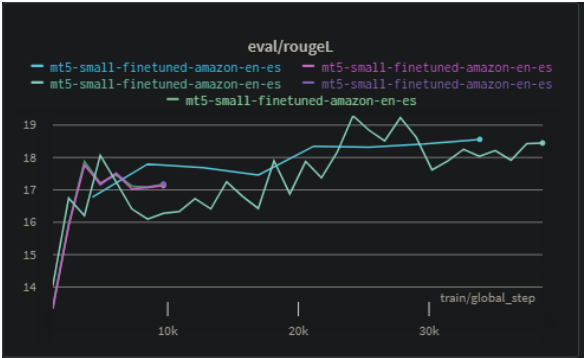


Figure 3: Rogue-L graph for runs [4]

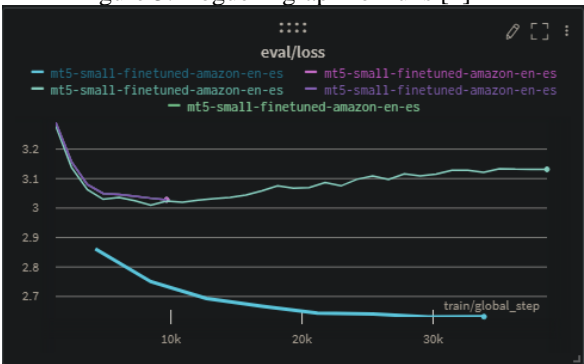


Figure 4: Loss graph [4]

181 I do not have saved data for the gradient
182 accumulation as that used a different model in the
183 code that both ended up being less accurate for
184 some reason and did not get logged with wandb. It
185 was putting out ROGUE-L 11 normally, with 9-10
186 with gradient accumulation.

187 **5 Conclusion**

189 Summarization is a task that takes large amounts of
190 GPU, data, and time to get proper results with.
191 Larger batches increase accuracy, as does longer
192 sequence length. Both of those take more GPU.
193 Working within the confines of the hardware I have
194 access to means fiddling with every
195 hyperparameter to see how that affects the results.
196 While there is not much more I can do with this
197 code, I can take the knowledge of optimization I
198 learned here and apply it to future machine learning
199 models.

References

Rush, Alexander & Chopra, Sumit & Weston, Jason. (2015). *A Neural Attention Model for Abstractive Sentence Summarization*. [1]

Yang, L. (2016). *Abstractive Summarization for Amazon Reviews*. [2]

Vaswani, Ashish & Shazeer, Noam & Parmar, Niki & Uszkoreit, Jakob & Jones, Llion & Gomez, Aidan & Kaiser, Lukasz & Polosukhin, Illia. (2017). *Attention Is All You Need*. [3]

Palmer, D. (n.d.). *Weights & Biases – developer tools for ML*. Weights & Biases – Developer tools for ML. Retrieved May 5, 2022, from <https://wandb.ai/site> [4]

Summarization - hugging face course. Summarization - Hugging Face. (n.d.). Retrieved May 5, 2022, from <https://huggingface.co/course/chapter7/5?fw=pt> [5]

Convolutional Neural Networks for Visual Recognition. CS231N convolutional neural networks for visual recognition. (n.d.). Retrieved May 5, 2022, from <https://cs231n.github.io/optimization-2/> [6]