

Methods

- Novel, “reverse engineering” approach to finding text summaries (find a better word... in an information sense, not in a writing short summary sense). Given a target sequence and a given LLM, we want to find the input embedding to this LLM, which maximizes the likelihood for the LLM to generate the target sequence.
- Training can be done with simply forward passes through the LLM model

For a given Large Language Model, \mathcal{M} , we want to find the input embedding, \mathbf{s}^* , that maximizes the likelihood of generating a given target sequence, $\{t_i\}_{i=1}^L$, where L denotes the length of the target sequence and t_i represents token i .

To express that we generate tokens with the LLM, based on length E input-embedding \mathbf{e} , we use the notation $\mathcal{M}(\mathbf{e})$. To express the number of tokens, k that we generate, we use the subscript $\mathcal{M}_k(\mathbf{e})$. This function call returns an output vector of length V , where V denotes the vocabulary size. We denote the output vector as \mathbf{o} . Each element o_v of \mathbf{o} represents the probability of generating token v as the next token. We predict the next token, by selecting the largest element of \mathbf{o} . Explicitly, $t_1 = \arg \max \mathcal{M}_1(\mathbf{e})$.

We denote the probability that \mathcal{M} generates the target sequence, given the input embedding, as $p(t_1, \dots, t_L | \mathbf{s})$.

We estimate these embeddings, by maximizing the conditional likelihood that this embedding generates the target sequence for a given LLM. We use a gradient-based optimization algorithm to maximize this likelihood.

These embeddings have the same dimension as the input embedding of \mathcal{M} , length E .

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} \prod_{i=1}^L p(t_i, \dots, t_L | \mathbf{s})$$

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} \sum_{i=1}^L \log p(t_i, \dots, t_L | \mathbf{s})$$

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} \log p(t_1, | \mathbf{s}) + \log p(t_2, | \mathbf{s}, t_1) + \dots + \log p(t_L, | \mathbf{s}, t_1, \dots, t_{L-1})$$

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} \log \mathcal{M}_1(\mathbf{s}) + \log \mathcal{M}_1([\mathbf{s}, \mathbf{e}_1]) + \dots + \log \mathcal{M}_1([\mathbf{s}, \mathbf{e}_1, \dots, \mathbf{e}_{L-1}])$$

$$\mathbf{s}^* = \arg \max_{\mathbf{s}} \sum_{i=1}^L \log \mathcal{M}_L(\mathbf{s})$$

We can make this expression more explicit, by using the LLM in our notation. For example $p(t_1, |\mathbf{s}) = \mathcal{M}(\mathbf{s})$

In the training process, we execute the following algorithm. First, we initialize the input embedding, \mathbf{s} . Second, we generate a sequence of length L , with the LLM and the current input embedding.

```

Initialize  $\mathbf{s}$ 
 $i \leftarrow 0$ 
while  $Convergence == False$  do
     $i \leftarrow i + 1$ 
     $\mathbf{o}^{(i)} \leftarrow \mathcal{M}(\mathbf{s})$ 
     $l^{(i)} \leftarrow \text{Likelihood}(\mathbf{o}^{(i)}, \{t_i\}_1^L)$ 
     $\mathbf{s}^{(i+1)} \leftarrow \text{GradientDecent}(\mathbf{s}^{(i)}, l^{(i)})$ 
     $Convergence \leftarrow \text{CheckConvergence}(\mathbf{s}^{(i+1)}, \mathbf{s}^{(i)})$ 
end while

```
