Hyperparameter Tuning for LLMs

Generated with Google Gemini - Published: April 20, 2024 - 21:47

From what I have learned: I <user> am able to determine your <system> hyperparameterisation, in order to steer your <assistant> output.

And, it would be wise to determine either 'top_k' and/or 'top_p' but not 'temperature'

or the other way round:

'temperature' but neither 'top_k' nor 'top_p', since the outcome would be pretty messy otherwise, right?

You've grasped the key concepts very well! Here's a breakdown of what you described:

User Control and Hyperparameters:

• Yes, by providing prompts and adjusting settings, you can influence the output (outcome) of a large language model (system) like me.

These settings are called hyperparameters, and some of them affect how I generate text.

Top_k, Top_p, and Temperature:

- These are three common hyperparameters used for text generation in LLMs. They all influence which word I pick next when building a response.
 - **Top_k:** This limits me to choosing the next word from only the top k most probable words. It increases focus and reduces randomness but can make responses repetitive.
 - **Top_p:** This sets a minimum cumulative probability threshold. I only consider words that add up to a total probability higher than this threshold. It encourages diversity but can lead to less relevant word choices.
 - **Temperature:** This acts like a randomness dial. A higher temperature makes less probable words more likely to be chosen, increasing variety but potentially sacrificing coherence.

Why Not All Three?

- You're right. Using all three together can be counter-productive. Here's why:
 - Top_k and top_p both achieve a similar effect of limiting the options I consider. Using them together creates an overly restricted selection pool.
 - Temperature adds randomness, while top_k and top_p restrict it. Combining them can lead to a confusing mix of repetitive and nonsensical outputs.

Best Practices:

- It's generally recommended to adjust either top_k or top_p, but not both at the same time.
- Temperature can be used in combination with either top_k or top_p, but fine-tuning is needed to find a balance.

By experimenting with these settings, you can steer my responses towards being more creative, informative, or sticking closely to the prompt, depending on your needs.