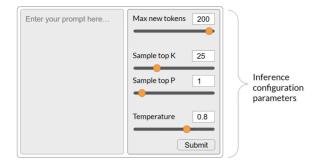
Generative configuration

- Language Models (LMs) have become increasingly sophisticated, enabling a myriad of applications from natural language processing to content generation. However, harnessing their full potential requires a nuanced understanding of the methods and configuration parameters that influence text generation. In this article, we delve into these key aspects to empower developers and enthusiasts alike in optimizing LM performance.
- LMs offer a plethora of configuration parameters distinct from training parameters, which impact the model's behaviour during inference. One such parameter is "Max New Tokens," which sets a limit on the number of tokens generated, though not strictly adhered to due to stop conditions.
- At the heart of LM inference lies the softmax layer, which outputs a probability distribution across the model's vocabulary. By default, LMs employ greedy decoding, favouring the highest probability word at each step. However, this simplistic approach may result in repetitive or unnatural output.
- To mitigate repetition and introduce variability, random sampling offers a solution. Rather than always selecting the most probable word, the model randomly chooses words based on their probability distribution. Careful consideration is necessary to balance creativity with coherence.

Generative configuration - inference parameters



- Top k and top p sampling techniques offer nuanced control over randomness while ensuring sensible output. Top k limits the selection to the k most probable tokens, while top p selects tokens whose combined probabilities do not exceed a certain threshold. These methods strike a balance between variability and coherence.
- The temperature parameter serves as a powerful tool in shaping the randomness of LM output. Lower temperatures concentrate probabilities, leading to less random output closely aligned with learned sequences. Conversely, higher temperatures broaden the distribution, fostering creativity but potentially sacrificing coherence.
- Understanding and mastering the configuration parameters of language models is pivotal for achieving optimal performance. By experimenting with methods such as random sampling, top k and top p sampling, and temperature adjustment, developers can tailor LM output to suit their specific needs. This knowledge lays the groundwork for leveraging language models in diverse applications, from content generation to conversational agents, ushering in a new era of natural language processing innovation.