# Top_k en Top_p : diversity and relevance ... what is this all about, basically?

What kind of hyperparameter is the Top_k value?

What kind of hyperparameter is the Top_p value?

Both Top_k and Top_p are hyperparameters used in controlling the randomness of text generation in large language models. They achieve this by influencing which word the model picks next when crafting a sequence of words. Here's a breakdown of each:

- **Top_k:** This is a discrete hyperparameter. It sets a fixed number (k) of the most likely words the model can choose from at each step during generation. Lower values of k restrict the model to a smaller set of high-probability words, resulting in less variation and more predictable outputs that closely align with the model's top predictions.

- **Top_p (nucleus sampling):** This is a continuous hyperparameter. It sets a probability threshold (p) for word selection. The model considers all the possible words and their probabilities. It then picks only those words whose cumulative probability reaches or surpasses the threshold (p). Finally, it randomly selects the next word from this narrowed-down set. Lower values of p lead to a smaller selection of top-weighted words, making the outputs less diverse and more focused on the model's most likely continuations.

Both hyperparameters trade-off between diversity and relevance in the generated text. They help you steer the model towards more creative or more conventional outputs depending on your needs.

That is what you say: Diversity and relevance ... ?!

Well, this seems something of a contradiction in NLP or rather in NLG, right?

A contradiction of an almost dialectic dimension, not true?

You're absolutely right. Diversity and relevance can indeed seem like opposing forces in Natural Language Generation (NLG). It's a fascinating challenge that researchers are constantly grappling with. Here's how we can think about this dialectic:

- **Diversity:** We want NLG models to be creative and explore a variety of possibilities when generating text. This can lead to surprising, interesting, and even humorous outputs. It injects richness and avoids monotonous repetition.

- **Relevance:** On the other hand, we also want the generated text to be coherent, on-topic, and aligned with the context. It should make sense and fulfill the purpose for which it was created.

**The crux of the problem lies in finding the right balance between these two forces.**

Here's how the Top_k and Top_p hyperparameters we discussed earlier play a role in this dialectic:

- **High Diversity, Low Relevance:** With a low Top_k or Top_p value, the model has more freedom to choose unexpected words, potentially leading to creative but less relevant outputs.

- **High Relevance, Low Diversity:** Conversely, high values for these hyperparameters restrict the model's choices, resulting in safe but potentially monotonous outputs that strictly adhere to the prompt or topic.

Finding the sweet spot depends on the specific application. Here are some scenarios:

- **Creative writing:** Here, diversity might be more important. A low Top_k or Top_p value could spark interesting ideas and unexpected turns of phrase.

- **Writing factual summaries:** In contrast, relevance is paramount. High values for these hyperparameters would ensure the generated text stays focused and accurate.

This dialectic extends beyond hyperparameter tuning. Researchers are exploring various techniques to improve NLG, such as:

- **Incorporating prior knowledge:** Injecting background information or specific domain knowledge into the model can help guide its word choices and enhance both diversity and relevance.

- **Refining objective functions:** The way the model is evaluated can influence its output. Metrics that consider both creativity and coherence can lead to a better balance.

The pursuit of diverse yet relevant NLG remains an ongoing area of research. It's a dialectic that pushes the boundaries of language models and opens doors to more sophisticated and nuanced forms of computer-generated text.