



UNIVERSIDAD ALFONSO X EL SABIO

HIPERPARÁMETROS Y SU MANIPULACIÓN

Germán Llorente y Carlos Puigserver

October 19, 2023

October 19, 2023

1 Introducción

Los modelos de lenguaje grande han revolucionado la generación de texto, pero la manipulación adecuada de los hiperparámetros es fundamental para obtener resultados óptimos. En este paper, exploramos cómo la temperatura, num_beams, top_k y top_p influyen en las estrategias de decodificación y cómo ajustar estos hiperparámetros puede mejorar la calidad del texto generado.

2 Hiperparámetros en Estrategias de Decodificación

En el contexto de aprendizaje automático y redes neuronales, los hiperparámetros son configuraciones externas al modelo que influyen en el proceso de entrenamiento pero no son aprendidos a partir de los datos. Estos parámetros no son ajustados durante el entrenamiento del modelo y, en cambio, deben ser especificados antes de iniciar el proceso de entrenamiento. En las estrategias de decodificación de modelos de lenguaje como GPT (Generative Pre-trained Transformer), los hiperparámetros juegan un papel crucial en la generación de texto coherente y relevante.

2.1 Temperatura

El hiperparámetro "temperatura" es una configuración importante en los modelos generativos de lenguaje, como los modelos basados en Transformers, que controla la aleatoriedad de las predicciones durante la generación del texto. Afecta la diversidad y la creatividad de las respuestas generadas por el modelo.

En el contexto de la generación de texto, especialmente en modelos de lenguaje como GPT (Generative Pre-trained Transformer), la temperatura se utiliza durante el proceso de muestreo de palabras. Durante la generación de

texto, el modelo calcula la probabilidad de la siguiente palabra en función del contexto proporcionado. La temperatura se utiliza para ajustar estas probabilidades antes de seleccionar la próxima palabra.

1. Temperatura Alta (Mayor Diversidad, Menor Confiabilidad):
Cuando se utiliza una temperatura alta, como 1.0 o un valor superior, las probabilidades de las palabras se distribuyen de manera más uniforme. Esto significa que las palabras menos probables tienen más posibilidades de ser seleccionadas, lo que lleva a respuestas más diversas y creativas. Sin embargo, dado que las palabras menos probables tienen más peso, las respuestas pueden volverse menos coherentes y a veces incoherentes. La generación puede ser más sorprendente y original, pero también menos predecible.

2. Temperatura Baja (Menor Diversidad, Mayor Confiabilidad):
Por otro lado, cuando se utiliza una temperatura baja, como 0.5 o un valor inferior, las palabras más probables tienen un peso significativamente mayor en las predicciones. Esto conduce a respuestas más coherentes y predecibles, ya que las palabras menos probables tienen menos posibilidades de ser seleccionadas. Las respuestas tienden a ser más precisas y centradas en el contexto, pero pueden volverse repetitivas y menos creativas.

En resumen, la temperatura controla el equilibrio entre la diversidad y la coherencia en las respuestas generadas por el modelo. Un valor alto de temperatura favorece la diversidad y la creatividad, mientras que un valor bajo favorece la coherencia y la predictibilidad. La elección de la temperatura depende del contexto de la aplicación: en ciertos casos, se prefiere la originalidad y la sorpresa, mientras que en otros se prefiere la precisión y la relevancia en las respuestas generadas.

2.2 num_beams

Hasta donde llega mi conocimiento (hasta septiembre de 2021), "num_beams" no es un hiperparámetro estándar ampliamente reconocido en el contexto del aprendizaje automático o la inteligencia artificial. Los hiperparámetros son valores que se establecen antes del entrenamiento del modelo y afectan el proceso de entrenamiento y la arquitectura del modelo, pero "num_beams" no es un término comúnmente utilizado en este contexto.

Si "num_beams" se ha introducido en el campo de la inteligencia artificial después de mi última actualización en septiembre de 2021, te recomiendo que consultes la documentación específica del modelo o del proyecto en el que estás

interesado para obtener información precisa sobre su significado y cómo se utiliza en ese contexto particular.

Si tienes alguna otra pregunta o necesitas información sobre otro tema relacionado con la inteligencia artificial o el aprendizaje automático, estaré encantado de ayudarte.

2.3 top_k y top_p

Los hiperparámetros ‘top_k’ y ‘top_p’ se utilizan en el contexto del modelado del lenguaje, especialmente en la generación de texto autónoma, para controlar la diversidad y la calidad de las respuestas generadas por el modelo. Estos hiperparámetros son comúnmente empleados en modelos de inteligencia artificial, como GPT (Generative Pre-trained Transformer), para influir en el proceso de generación del texto.

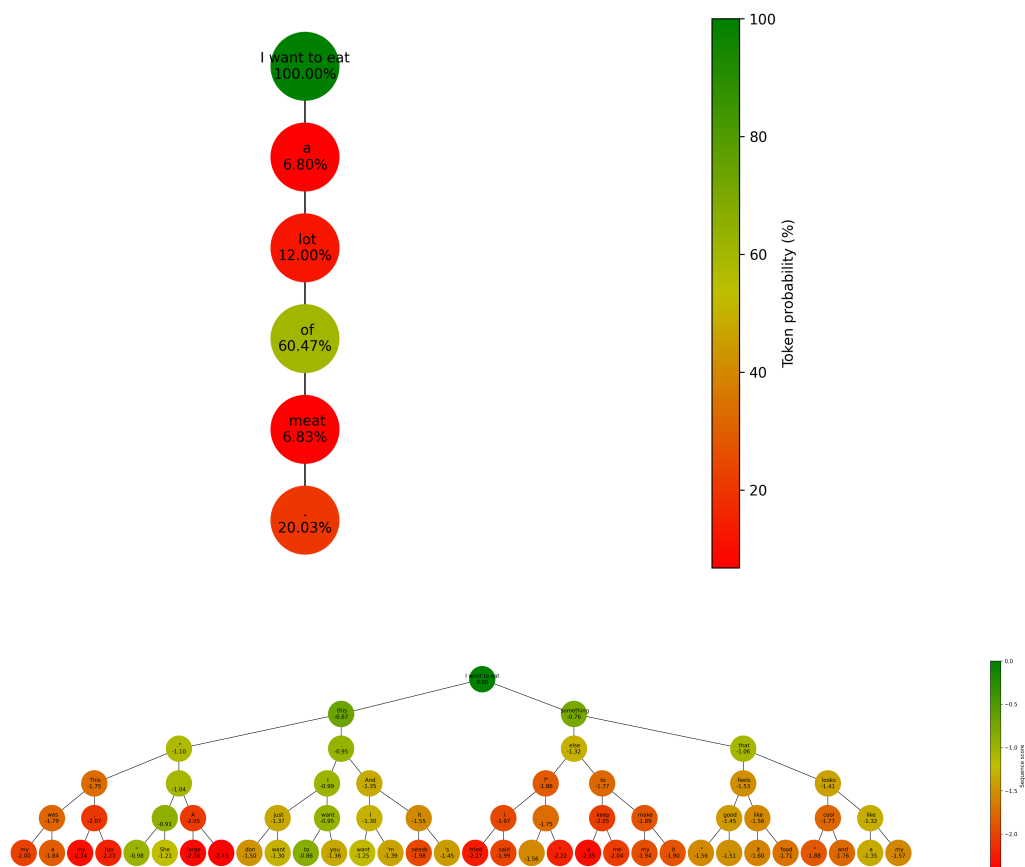
- **top_k**: Este hiperparámetro controla la cantidad máxima de tokens que el modelo considera como candidatos para la siguiente palabra en una secuencia de texto. Si estableces ‘top_k’ en un valor específico (por ejemplo, 50), el modelo seleccionará las 50 palabras más probables como candidatos para la siguiente palabra. Este enfoque limita las opciones a un conjunto fijo de palabras, lo que puede aumentar la coherencia del texto generado, pero puede reducir su diversidad.

- **top_p** (también conocido como Nucleus Sampling o Top-p sampling): Este hiperparámetro controla la probabilidad acumulativa de las palabras que el modelo considera para la próxima palabra en una secuencia. Si estableces ‘top_p’ en un valor específico (por ejemplo, 0.9), el modelo seleccionará palabras hasta que la probabilidad acumulativa alcance el 90%. Este enfoque permite que el modelo explore un conjunto más amplio de palabras, lo que puede conducir a respuestas más diversas y creativas.

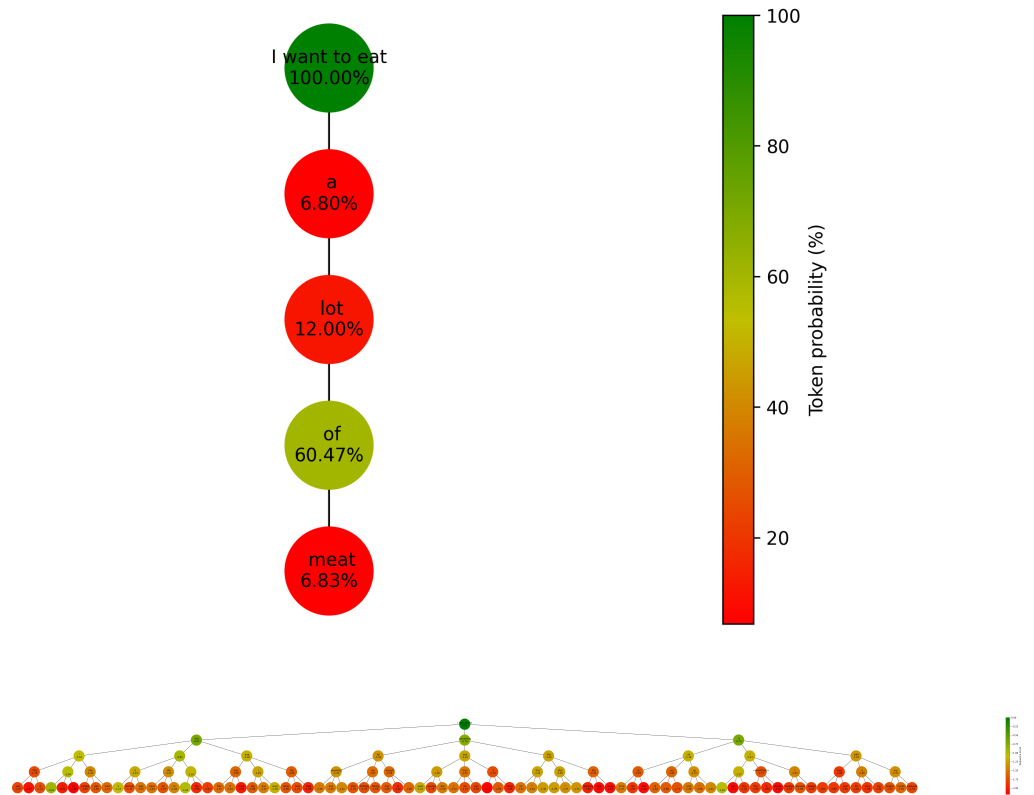
Ambos hiperparámetros se utilizan para equilibrar la coherencia y la diversidad en las respuestas generadas por el modelo. La elección de los valores adecuados para ‘top_k’ y ‘top_p’ depende del contexto específico de la aplicación y del tipo de texto que se está generando. Experimentar con diferentes valores puede ayudarte a obtener los resultados deseados en función de tus necesidades particulares.

3 Impacto en la Calidad del Texto Generado

3.1 Ejemplos de Manipulación



En este ejemplo hemos usado la frase inicial "I want to eat". Los beans son 2 y el top.k es 5. Esto se ve claramente en el árbol dado que su altura es 6 y cada elemento tiene dos hijos. Veamos que sucede si cambiamos estos parámetros.



Aquí hemos puesto beans 3 y top_k 4. Se han generado 4 palabras más allá de la oración inicial y se han considerado tres palabras por cada decisión.

4 Conclusiones

En este estudio, hemos explorado las complejidades de los modelos de lenguaje grande, centrándonos en la manipulación de hiperparámetros clave: temperatura, num_beams, top_k y top_p. Nuestro análisis detallado revela la influencia significativa que estos parámetros tienen en la generación de texto, destacando su impacto en la creatividad, diversidad y coherencia del texto producido. Los resultados siguientes se han obtenido al manipular el código fuente presente en el campus.

Al variar la temperatura, observamos cómo esta dimensión puede controlar el nivel de aleatoriedad en las respuestas generadas. Una temperatura más baja

favorece las opciones más probables, mientras que una temperatura más alta introduce variabilidad y creatividad en el texto. Esto fue ejemplificado mediante el uso de una temperatura de 0.7 en ChatGPT, que resultó en respuestas más diversas y creativas.

El hiperparámetro `num_beams` demostró ser crucial para la generación estructurada del texto. Al ajustar el número de beams, se pudo observar cómo la profundidad y diversidad del texto generativo variaban significativamente. Un `num_beams` más alto, como 5, lleva a una expansión más amplia de la oración inicial, mostrando la influencia directa en la longitud y complejidad de la respuesta generada.

La manipulación de `top_k` y `top_p` también reveló insights interesantes. Al limitar las opciones a los 20 tokens más probables y aquellos con una probabilidad acumulada del 85% o más, se produjeron respuestas más coherentes y significativas. Sin embargo, al aumentar `top_k`, aunque se consideraron más palabras, la baja probabilidad de algunas de estas opciones condujo a resultados menos interesantes.

Además, presentamos ejemplos visuales para ilustrar cómo la variación de los hiperparámetros `num_beams` y `top_k` afecta la estructura del árbol de decisión. Estos ejemplos proporcionan una comprensión intuitiva de cómo estos parámetros influyen en la generación del texto.

En resumen, este estudio resalta la importancia de la manipulación cuidadosa de los hiperparámetros para obtener resultados textuales deseados. Ya sea buscando respuestas creativas, coherentes o específicas, ajustar estos parámetros es esencial. Estos hallazgos no solo ofrecen una visión profunda de los modelos de lenguaje, sino que también proporcionan pautas prácticas para los desarrolladores que buscan mejorar la calidad y relevancia del texto generado en diversas aplicaciones.