

## Codificación posicional en los Transformers visuales: Capturando el contexto espacial en el análisis de imágenes

Los Transformers visuales han revolucionado el campo del procesamiento de imágenes y el análisis visual, logrando avances significativos en tareas como la clasificación de imágenes, la detección de objetos y la segmentación. Estos modelos, basados en la arquitectura del Transformer, han demostrado su capacidad para capturar relaciones contextuales en las imágenes, permitiendo un procesamiento más eficiente y efectivo de la información visual.

Sin embargo, a pesar de su éxito, los Transformers visuales presentan una limitación inherente: carecen de información intrínseca sobre la ubicación espacial de los elementos en las imágenes. A diferencia de los enfoques tradicionales basados en convoluciones, que aprovechan la estructura espacial de la imagen, los Transformers tratan las imágenes como secuencias planas de parches, perdiendo así la noción espacial en el proceso.

La captura precisa de la ubicación espacial es esencial para el análisis de imágenes, ya que a menudo la información contextual depende de la posición relativa de los objetos y las características visuales. Para abordar esta limitación, se ha desarrollado una técnica clave en los Transformers, que luego fue adaptada para los Transformers visuales: la codificación posicional.

En esta lectura, exploraremos en detalle la codificación posicional en los Transformers visuales y cómo se utiliza para capturar la información espacial en el análisis de imágenes. Analizaremos las estrategias comunes utilizadas en la codificación posicional, así como sus ventajas y desafíos.

A medida que profundicemos en los conceptos de codificación posicional, descubriremos cómo esta técnica se convierte en un componente esencial para aprovechar al máximo los Transformers visuales en el análisis de imágenes. Al final de esta lectura, esperamos que se pueda comprender mejor cómo la codificación posicional enriquece la capacidad de los Transformers visuales para modelar relaciones espaciales y contextualizar la información visual.

## Limitaciones de los Transformers visuales en la captura de información posicional

A pesar de su éxito en el análisis de imágenes, los Transformers visuales presentan una limitación inherente en la captura de información posicional. Esto se debe a que los Transformers se diseñaron originalmente para procesar secuencias de texto, donde el orden de las palabras es fundamental. Al aplicarlos al análisis de imágenes, se enfrentan al desafío de capturar información sobre la ubicación espacial de los objetos en una imagen.

Esta limitación se vuelve especialmente relevante en tareas donde la información contextual depende de la posición relativa de los objetos. Por ejemplo, en la detección de objetos, saber la ubicación precisa de cada objeto es fundamental para realizar una detección precisa y evitar falsos positivos. Además, en la segmentación semántica, la comprensión de la relación espacial entre diferentes regiones de la imagen es esencial para una segmentación precisa.

Sin una captura precisa de la información posicional, los Transformers visuales podrían enfrentar dificultades para modelar adecuadamente las relaciones espaciales y contextualizar la información visual. Por lo tanto, es crucial introducir estrategias de codificación posicional en los Transformers visuales para superar esta limitación y mejorar su capacidad para comprender la estructura espacial de las imágenes.

## Codificación posicional en Transformers visuales

La codificación posicional es una técnica fundamental en los Transformers, que luego fue adaptada al contexto de los Transformers visuales para capturar información sobre la ubicación espacial de los elementos en una imagen. Esta estrategia permite enriquecer la representación de cada elemento con información de posición, lo que mejora la capacidad del modelo para comprender las relaciones espaciales y contextualizar la información visual.

En los Transformers visuales, una estrategia comúnmente utilizada para codificar la información espacial es mediante el uso de embeddings de posición, también llamados embebimientos o codificadores posicionales. Los embeddings de posición son vectores de números reales que

representan la posición relativa de cada parche en la imagen. Estos embeddings se combinan con las características visuales de cada parche y se introducen en el modelo Transformer.

La generación de los embeddings de posición puede realizarse de diferentes maneras. Una opción es utilizar funciones trigonométricas como seno y coseno para asignar valores a cada dimensión del embedding. Por ejemplo, se puede asignar el valor del seno a una dimensión del embedding para codificar la posición horizontal, y el valor del coseno a otra dimensión para codificar la posición vertical. Esta aproximación es similar a los embeddings de posición en los Transformers para procesamiento de lenguaje natural.

Otra opción es aprender los embeddings de posición como parte del entrenamiento del modelo. Durante el proceso de aprendizaje, el modelo ajusta los parámetros de los embeddings de posición para capturar las relaciones espaciales relevantes en los datos de entrenamiento. Esto permite que el modelo adquiera una representación más precisa de la información posicional y mejore su capacidad para comprender la estructura espacial de las imágenes.

Una vez generados los embeddings de posición, se concatenan con las características visuales de cada parche antes de pasarlos al modelo Transformer. Esta concatenación enriquece la representación de cada parche con información de posición, lo que permite al modelo capturar patrones relacionados con la ubicación espacial.

La introducción de los embeddings de posición en los Transformers visuales ayuda al modelo a distinguir y comprender la información espacial en las imágenes. Al capturar la información de posición, el modelo puede reconocer la ubicación relativa de los objetos y comprender mejor las relaciones espaciales y contextuales.

## Ventajas y desafíos de la codificación posicional

La codificación posicional en los Transformers visuales ofrece varias ventajas significativas en el análisis de imágenes, pero también presenta desafíos que deben abordarse. A continuación, exploraremos tanto las ventajas como los desafíos asociados con la codificación posicional en los Transformers visuales.

### **Ventajas de la codificación posicional:**

1. **Captura de relaciones espaciales:** La codificación posicional permite que los Transformers visuales capturen de manera más efectiva las relaciones espaciales entre los elementos en una imagen. Al incorporar información sobre la ubicación relativa de los objetos, el modelo puede comprender mejor la estructura espacial de la imagen y utilizar esta información para tareas como detección de objetos, segmentación semántica y generación de descripciones.
2. **Contextualización espacial:** La información posicional codificada en los Transformers visuales ayuda a contextualizar las características visuales de cada elemento. Esto permite que el modelo comprenda el contexto espacial en el que se encuentran los objetos, lo que es crucial para comprender el significado y la relación entre diferentes partes de la imagen.
3. **Generalización a diferentes tamaños y aspectos:** Al codificar la información posicional, los Transformers visuales adquieren una mayor capacidad de generalización a diferentes tamaños y aspectos de las imágenes. Esto significa que el modelo puede manejar imágenes de diferentes resoluciones y relaciones de aspecto sin comprometer su capacidad para capturar relaciones espaciales y comprender la estructura de la imagen.

### **Desafíos de la codificación posicional:**

1. **Dependencia de la resolución y escala:** La codificación posicional puede verse afectada por la resolución y escala de la imagen. Los Transformers visuales pueden enfrentar dificultades al capturar relaciones espaciales finas en imágenes de baja resolución o con objetos de diferentes tamaños. Se requiere un diseño cuidadoso de los embeddings de posición y una elección adecuada de la escala para abordar este desafío.
2. **Sensibilidad a la ubicación precisa:** La codificación posicional puede ser sensible a pequeñas variaciones en la ubicación precisa de los objetos en la imagen. Dado que los Transformers visuales tratan los objetos como parches independientes, pueden enfrentar dificultades al lidiar con objetos superpuestos o parcialmente ocultos. La precisión en la

localización de los objetos puede ser crucial para el rendimiento óptimo en tareas que requieren una comprensión espacial precisa.

3. Combinación con otras técnicas: La codificación posicional en los Transformers visuales debe combinarse adecuadamente con otras técnicas y enfoques para obtener el mejor rendimiento en tareas específicas. Es importante considerar cómo interactúa la codificación posicional con las características visuales extraídas, las operaciones de atención y otros componentes del modelo.

## Bibliografía

- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
- Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... & Guo, B. (2021). Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF international conference on computer vision (pp. 10012-10022).



---

© - **Derechos Reservados:** la presente obra, y en general todos sus contenidos, se encuentran protegidos por las normas internacionales y nacionales vigentes sobre propiedad Intelectual, por lo tanto su utilización parcial o total, reproducción, comunicación pública, transformación, distribución, alquiler, préstamo público e importación, total o parcial, en todo o en parte, en formato impreso o digital y en cualquier formato conocido o por conocer, se encuentran prohibidos, y solo serán lícitos en la medida en que se cuente con la autorización previa y expresa por escrito de la Universidad de los Andes.

De igual manera, la utilización de la imagen de las personas, docentes o estudiantes, sin su previa autorización está expresamente prohibida. En caso de incumplirse con lo mencionado, se procederá de conformidad con los reglamentos y políticas de la universidad, sin perjuicio de las demás acciones legales aplicables.

---