

Les transformers sont des systèmes qui ont révolutionné le monde de l'apprentissage de langage notamment en 2017 où Google a sorti un article pour dire que cela était plus efficace que tout les autres par la suite Google a sorti son LLM puis OpenAI en 2018.

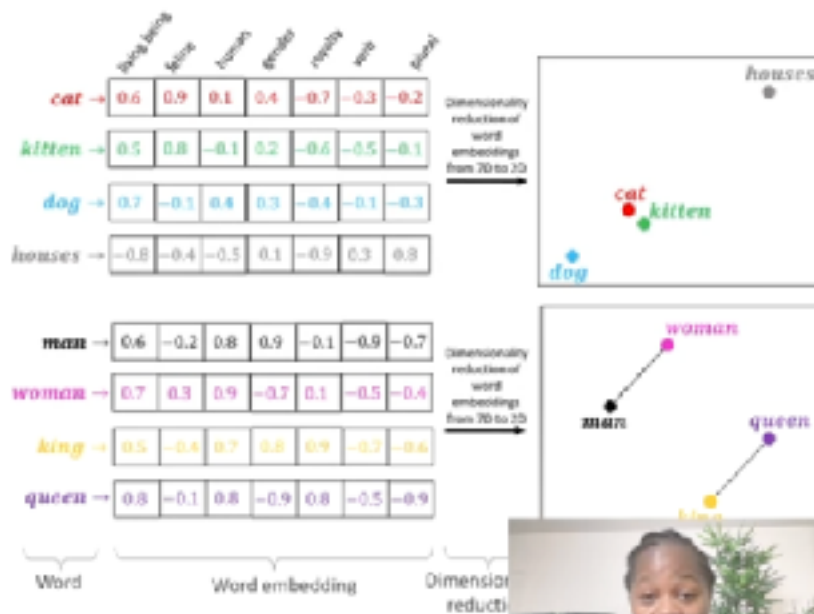
Tokenisation

Cela est un processus de découpage d'un texte en morceaux plus petits comme des mots des lettres ou des phrases.

En deep learning il faut être capable de transformer ces mots en représentation vectorielle cela s'appelle l'Embedding

Embedding

L'embedding est une représentation vectorielle qui capture la sémantique des mots ou des phrases dans un espace de haute dimension



Comme ici chaque mot aura une représentation vectorielle en 7D qui sera ensuite représentée en 2D comme sur le graphique

Chaque dimension représente une catégorie par exemple la partie humaine le genre etc ...

Cela va pouvoir permettre de manipuler ces mots en faisant des calculs dessus etc ...

NLP

En deep learning il y a 2 parties l'encodage et le décodage

Encodeur

L'encodeur est l'input c'est à dire qu'il va récupérer ce qu'il y a en entrée pour le transformer et pouvoir le projeter dans un espace par exemple transformé la phrase en la séparant en mot puis transformé les mots en vecteurs pour pouvoir ensuite les manipuler.

Ca peut être de la convolution CNN comme vue dans l'algorithme de convolution des images cela permettait de transformer les images avec différents filtres donc ça crée par exemple 4 images avec 4 images différentes qui mettent en avant par exemple les bordures ou les textes ... ça peut également modifier le padding ou le cadrage de l'image. Puis ensuite il y a l'effet filtre qui efface les pixels négatifs les transforme en 0 etc ... Puis enfin le dernier filtre qui permet de réduire la taille de l'image en combinant des pixels ensemble ça permet de réduire par 2 les pixels et d'optimiser le temps d'apprentissage.

Décodeur

Le décodeur son rôle est de récupérer tout ça pour faire la sortie en utilisant un modèle pour entraîner l'IA

Transformers

Les transformers par rapport aux anciennes techniques est beaucoup plus efficace notamment grâce à la parallélisation. Il va par exemple traiter la phrase en entier et non traiter chaque mot 1 par 1.

Cela se fait grâce à différents mécanismes.

Mécanisme d'attention

Le mécanisme d'attention est une technique permettant aux réseaux de neurones de se concentrer sur une partie spécifique de la séquence d'entrée. C'est à dire qu'elle va donner plus d'attention à certains mots de la phrase grâce à une formule traduite sous python.

Multi-head Attention

Une méthode qui utilise plusieurs "têtes" d'attention en parallèle pour capturer diverses relations dans les données.

Cela va donc diviser les vecteurs par le nombre de tête puis refaire le mécanisme d'attention c'est comme un système parallèle au mécanisme d'attention.