# Transformers

| Plan | Logistics |
|---|---|
| Review | check in |
| Self-Attention | scribe |
| Cross-Attention | Zoom |
| Positional Encoding | |

## Review

Motivation: words as vectors



one-hot encoding

$$f: \mathbb{R}^{|V|} \to \mathbb{R}^d$$

$(x, x^+)$ close $\Rightarrow f(x) \cdot f(x^+)$ large

$(x, x^-)$ far $\Rightarrow |f(x) \cdot f(x^-)|$ small

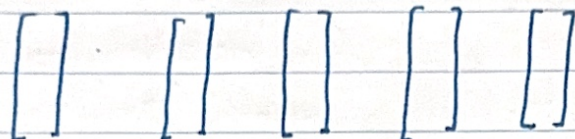C Contrastive learning $\subseteq$ unsupervised learning

Principal Component Analysis



capture variation of ~~eigenvectors~~ data

via eigenvectors

$$X^T X = \sum_{i=1}^{r} \lambda_i \, v^{(i)} v^{(i) \, T} \qquad \text{for} \quad v^{(i)} \cdot v^{(j)} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$

$$\max_{v: \|v\|_2 = 1} \quad \|Xv\|_2^2 \leftarrow v^T X^T X v \leftarrow \sum_{i=1}^{n} \lambda_i [v^T v^{(i)}]^2$$

Motivation: sentences as vectors

"Vermont is chilly and beautiful"

$$\begin{bmatrix} \\ \\ \end{bmatrix} \quad \begin{bmatrix} | \\ | \end{bmatrix} \quad \begin{bmatrix} \\ \\ \end{bmatrix} \quad \begin{bmatrix} \\ \\ \end{bmatrix} \quad \begin{bmatrix} \\ \end{bmatrix}$$

How can we understand __sequences__ of vectors?
↳ Recurrent networks
↳ LSTM

Attention! (self first)    $X \in \mathbb{R}^{n \times d}$

Goal: Combine similar words/tokens

Queries: $W^{(Q)} X = Q$    $W^{(Q)} \in \mathbb{R}^{r \times n}$
Keys: $W^{(K)} X = K$
Values: $W^{(V)} X = V$    $W^{(V)} \in \mathbb{R}^{d \times n}$

attention : $\text{softmax}(Q K^T) = \begin{bmatrix} \\ \\ \\ \end{bmatrix}$ — $Q_i^T K_j$
     ↑
   to rows

     $i \rightarrow$     $n \times n$

result: $\text{softmax}(Q K^T) V$

$\sum\limits_{i=1}^{n} \text{sim}(j,i) \, v_i$

$j \begin{bmatrix} \\ \\ \\ \end{bmatrix} \begin{bmatrix} \\ \\ \\ \end{bmatrix} = j \begin{bmatrix} \\ \\ \\ \end{bmatrix}$

$n \times n$       $n \times d$

$X \begin{bmatrix} \\ \\ \\ \end{bmatrix}$ → $\begin{cases} Q \begin{bmatrix} & \end{bmatrix} \\ K^T \begin{bmatrix} \\ \end{bmatrix} \\ V \begin{bmatrix} & \end{bmatrix} \end{cases}$ $n \times n$   $n \times d$  softmax  Attention $\begin{bmatrix} \\ \end{bmatrix}$ →

# Cross - Attention!

"Vermont is chilly and beautiful"

$$[\ ] \quad [\ ] \quad [\ ] \quad [\ ] \quad [\ ] \qquad X \begin{bmatrix} \\ \\ \end{bmatrix} \; n \times d$$

"Vermont es fria y"

$$[\ ] \quad [\ ] \quad [\ ] \quad [\ ] \qquad Y \begin{bmatrix} \\ \\ \end{bmatrix} \; m \times d$$

Goal: Represent sequence as linear combo of another

Queries:
$$W^{(Q)} X = Q \qquad\qquad W^{(Q)} \in \mathbb{R}^{r \times m}$$
$$W^{(k)} X = K \qquad\qquad W^{(k)} \in \mathbb{R}^{r \times n}$$
$$W^{(v)} X = V \qquad\qquad W^{(v)} \in \mathbb{R}^{d \times n}$$



$$\text{softmax}(K^T Q)$$

$$\sum_{i=1}^{n} \text{sim}(j, i)\, V_i$$

$m \times d$

# Large Language Models



Make upwards direction?

"Hello! What am ?!"

CLS ... SEP

self attention

FC  FC  FC

FC

↑ distribution over next words

## Positional Encoding?

We represent time as
11:24am Tuesday, Jan 14, 2025
rather than
1,065,066,880 min since 0 BC

↳ min captures schedule
↳ hour captures time of day
↳ day captures schedule
↳ date captures schedule
↳ month captures time of year
↳ year captures years passed

"Hello! What am ?!"
0    1    2    = t

$\sin(8^3 t)$
$\sin(6^2 t)$
$\sin(a_5 t)$