# Transformers
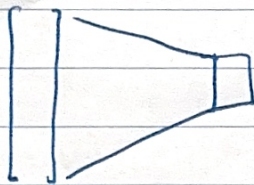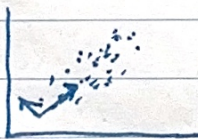
## Review

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



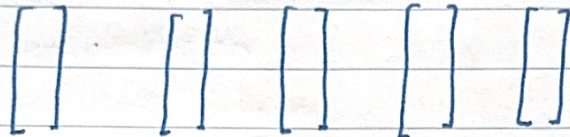capture variation of data eigenvectors

via eigenvectors

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

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

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: $W^{(Q)} X = Q$        $W^{(Q)} \in \mathbb{R}^{r \times d}$

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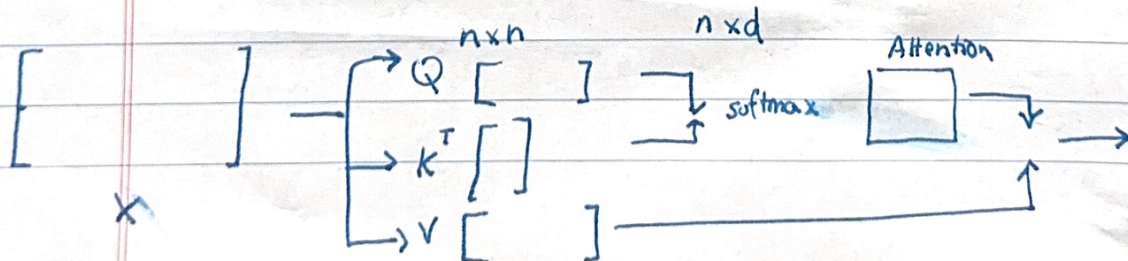
Values: $W^{(V)} X = V$        $W^{(V)} \in \mathbb{R}^{d \times d}$

attention : $\text{softmax}\left( Q K^T \right) =$

to rows

$i \rightarrow$

$Q_i^T K_j$

$n \times n$

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

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

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

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

Attention

$X \begin{bmatrix} \\ \\ \end{bmatrix} \longrightarrow \begin{cases} Q \begin{bmatrix} \\ \end{bmatrix} \\ K^T \begin{bmatrix} \\ \end{bmatrix} \\ V \begin{bmatrix} \\ \end{bmatrix} \end{cases}$  softmax  $\longrightarrow \boxed{\phantom{xx}} \longrightarrow$

# Cross - Attention!

"Vermont is chilly and beautiful"

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

"Vermont es fria y"

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

Goal: Represent sequence as linear combo of another

$$X \in \mathbb{R}^{n \times d}$$

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



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

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

$$m \times d$$

# Large Language Models

"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

$\sin(2^0 t)$

$\sin(2^1 t)$

$\sin(2^2 t)$

t