# 15-780: Graduate AI Lecture 9: transformers

Aditi Raghunathan

# Logistics

- Homework 2
  - Grading complete
  - Solutions out this evening
- Homework 3
  - Due next Monday
- Midterm
  - Everything including todays and upcoming Wed lecture
  - Next Monday: **review session (optional attendance)**

# Recap of deep networks

$\rightarrow$ MLP $\qquad z_i = \sigma(w_i^T z_{i-1})$

$\Rightarrow$ expressivity

$\rightarrow$ running gradient descent

auto diff framework $\qquad$ chain rule

# Batching in deep learning

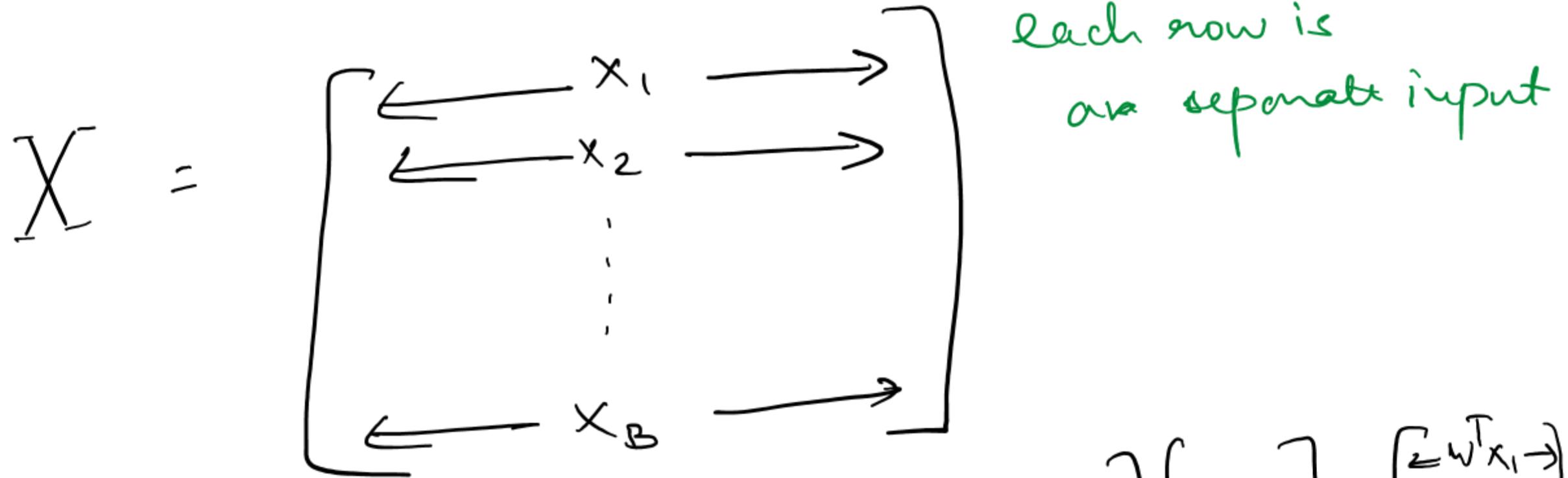Recap of stochastic gradient descent

Have a batch of examples $B$ $\rightarrow$ it entire train set = gradient descent

In each step

$$w_i^{(t)} = w_i^{(t-1)} - \sum \eta \nabla loss(x, y; w)$$
$$x, y \in B$$

key idea: $\nabla loss(x_i, y_i; w)$ is independent of other examples

$x_i \in \mathbb{R}^d$  B examples in a batch

$$X = \begin{bmatrix} \longleftarrow & x_1 & \longrightarrow \\ \longleftarrow & x_2 & \longrightarrow \\ & \vdots & \\ \longleftarrow & x_B & \longrightarrow \end{bmatrix}$$

<span style="color:green">each row is are separate input</span>

$W^T x \Longrightarrow X W$

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

$$\begin{bmatrix} \longleftarrow & x_1 & \longrightarrow \\ \longleftarrow & x_2 & \longrightarrow \\ & \vdots & \\ \longleftarrow & x_B & \longrightarrow \end{bmatrix} \begin{bmatrix} W \end{bmatrix} = \begin{bmatrix} \longleftarrow & W^T x_1 & \longrightarrow \\ \longleftarrow & W^T x_2 & \longrightarrow \\ & \vdots & \\ \longleftarrow & W^T x_B & \longrightarrow \end{bmatrix}$$

$$\underline{X} W \quad : \quad \underline{X} = \begin{bmatrix} \longleftarrow x_1 \longrightarrow \\ \longleftarrow x_2 \longrightarrow \\ \vdots \\ \longleftarrow x_B \longrightarrow \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} a' & b' \\ c' & d' \end{bmatrix} = \begin{bmatrix} aa' + bc' & ab' + bd' \\ ca' + dc' & cb' + dd' \end{bmatrix}$$

# Structured inputs

- Images

28

28

each pixel & roll out

784

$[ x, y \quad - \quad - \quad - \quad - \quad - \quad ]$

pixels

⇓ missing out
spatial information

- Text

*The quick brown fox jumps* over the lazy dog

# Key idea

We need to reason about **sets** of inputs

→ collection of pixels
→ collection of words

# Language modeling

$$e_{word} \in \{0,1\}^V$$

- Notation for one-hot vector

$e_{word}$ = zero every where except in the location of word in vocabulary

- Vocabulary $V$

- Input: **sequence** of T tokens

words

$e_{the}$, $e_{quick}$, $e_{brown}$

$$X = \begin{bmatrix} \longleftarrow e^T_{word_1} \longrightarrow \\ \longleftarrow e^T_{word_2} \longrightarrow \\ \vdots \\ \longleftarrow e^T_{word_T} \longrightarrow \end{bmatrix} \qquad X \in \mathbb{R}^{T \times V}$$

we want $p(word_{T+1} \mid X) = \mathbb{R}^V$

probability distribution over $V$ words

# Batch operations?

$$X W$$

$$
\begin{bmatrix}
\leftarrow \ell\ \text{the} \longrightarrow \\
\leftarrow \ell\ \text{quick} \longrightarrow \\
\vdots
\end{bmatrix}
W =
\begin{bmatrix}
\leftarrow \ell\ \text{the}^T\ w \longrightarrow \\
\leftarrow \ell\ \text{quick}^T\ w \longrightarrow \\
\vdots
\end{bmatrix}
$$

# Piazza poll

$$X = \begin{bmatrix} \longleftarrow \ell \text{ word}_1 \longrightarrow \\ \longleftarrow \ell \text{ word}_2 \longrightarrow \\ \vdots \end{bmatrix} \Big\}$$ same example

- Which of the following operations allow for sharing information across words

$A, W$ are some matrices of app dim

$\longrightarrow$ batching

- (A) $XW$
- (B) $\sigma(XW)$
- (C) $A\,X$ $\longrightarrow$
- (D) $\sigma(AX)$

$A(X) \; X$

A matrix depends on $X$ for self-attention

# Mixing information

$$\begin{bmatrix} a' & b' \\ c' & d' \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a'a + b'c & a'b + b'd \\ c'a + d'c & c'b + d'd \end{bmatrix}$$

$AX$ : combines info across rows

$XW$ : treats each row separately

impractical

$$\begin{bmatrix} \leftarrow \text{word}_1 \rightarrow ; & \leftarrow \text{word}_2 \rightarrow & - - - - \end{bmatrix}$$

# A: a probability distribution

$$Y = A \qquad X$$

$$X \in \mathbb{R}^{T \times d}$$
$$A \in \mathbb{R}^{T \times \boxed{T}}$$
$$Y \in \mathbb{R}^{T \times d}$$

$$\left( \begin{bmatrix} \alpha_1 & \alpha_2 \\ \beta_1 & \beta_2 \end{bmatrix} \begin{bmatrix} \leftarrow \ell word_1 \rightarrow \\ \leftarrow \ell word_2 \rightarrow \end{bmatrix} \right.$$

$$\begin{bmatrix} \alpha & 1-\alpha \\ \beta & 1-\beta \end{bmatrix} \begin{bmatrix} \leftarrow \ell wod 1 \rightarrow \\ \leftarrow \ell word2 \rightarrow \end{bmatrix} = \begin{bmatrix} \alpha \ell word 1 + (1-\alpha) \ell word_2 \\ \beta \ell word1 + (1-\beta) \ell word2 \end{bmatrix}$$

$P_i$'s are
dist over
$T$ elements

$$\begin{bmatrix} \leftarrow P_1 \rightarrow \\ \leftarrow P_2 \rightarrow \\ \vdots \\ \leftarrow P_T \rightarrow \end{bmatrix} \begin{bmatrix} \leftarrow 1 \rightarrow \\ \vdots \\ \leftarrow T \rightarrow \end{bmatrix} = \quad T \text{ diff combinations} \\ \text{of } X$$

# A: a probability distribution

$$Y = A X$$

$$\downarrow$$

- Each row is a probability distribution

- Each row $i$ of $X$ corresponds to $i^{th}$ word

# Creating the matrix A

- construct scores $\longrightarrow$ softmax

- " similarity or dot product
   between words $i$ & $j$ "

$$P: i^{(i)}_{\text{row}} \left[ \leftarrow \quad i \quad \rightarrow \right] \left[ \qquad \right]$$

$\sum P_j^{(i)} \; \ell \text{ word}_j$

$P_j^{(i)}$ : similarity
b/w $i$ & $j$

# Creating the matrix A

$$A = \text{softmax}\left( \frac{(X W_K) \cdot (X W_Q)^T}{\sqrt{d}} \right)$$

$$P_j^{(i)} = \left( X_i W_Q \right)^T \left( X_j W_K \right)$$

query       key

# Final form for self-attention

$$Y = \text{Softmax}\left( \frac{X \, W_Q \, W_K^T \, X^T}{\sqrt{d}} \right) X \, W_V$$

$W_A$

For now: $W_Q$, $W_K$ & $W_V$ are $\mathbb{R}^{d \times d}$ matrix

# Properties of attention

- Full mixing

$$p_j^{(i)} : \text{looks at } \text{``}word_i\text{,'' \& ``}word\,j\text{''})$$

$$\underline{\vee} \; i, j \; \in [1 \cdots T]$$

# Properties of attention

- Let us increase the size of the set T
  - What happens to $W_q, W_k, W_v$ ?   $\leftarrow \mathbb{R}^{d \times d}$

• $\mathbb{R}^{T \times T}$   $O(T^2)$ in the construction of $A$

# Properties of attention

- Does the ordering between words matter?

Fixed $W_Q, W_{K}, W_V$

does $A$:

"The quick brown fox" : $x$

" the brown fox quick" : $\tilde{x}$

$A(x)$ & $A(\tilde{x})$ are same upto permutations