

## Week 7

10/7/2025

↳ Quiz

↳ Last Week: 26, 33

↳ Midterm 10/21

- practice exam(s) coming soon
- Linear Algebra → Transformers
- 6 multiple choice, three long

Context:

Supervised Learning  
Models {

- Linear Regression
- Logistic Regression
- Support Vector Machines

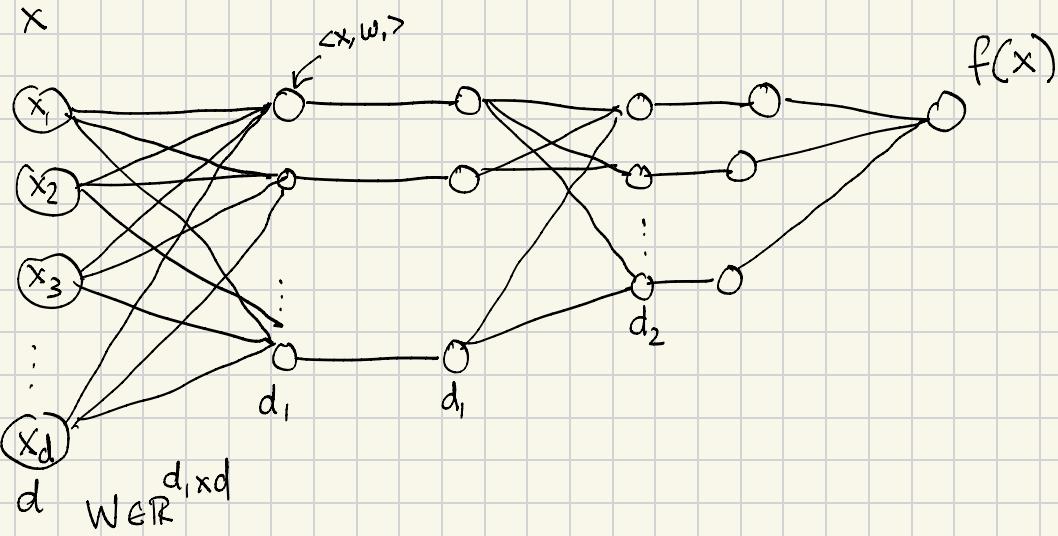
Feature Selection {

- Feature Transformation
- Kernel Trick
- Reparameterization Trick

=> Neural Networks!

- Convolutional (images, audio, etc)
- Transformers (text, sequential)

## Review : Fully Connected NNs



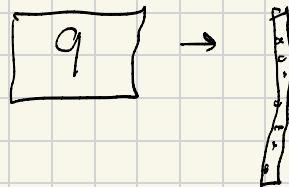
`nn.Linear(d, di)    nn.ReLU() ...`

$W_x = z$      $\sigma(z)$     ...

## Issues with Linear Layers

1. Computationally expensive

2. Loss of "context"

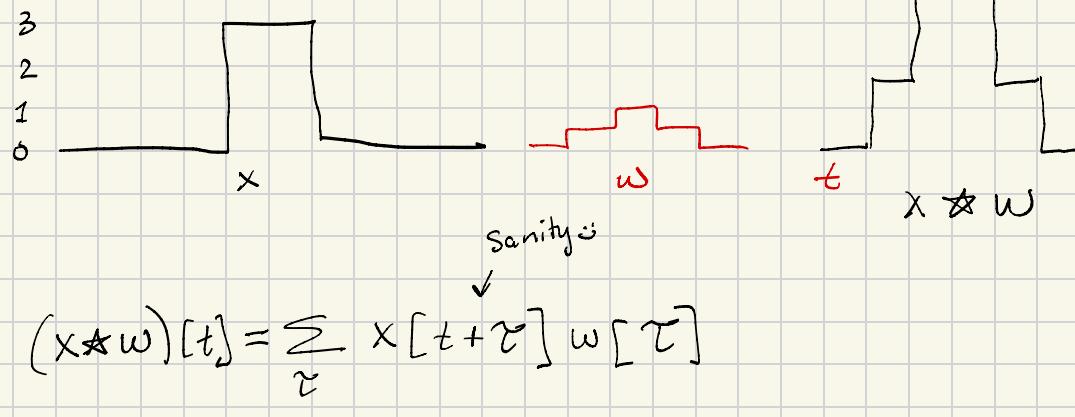


## Wishlist for Convolutional Layers

1. Locality

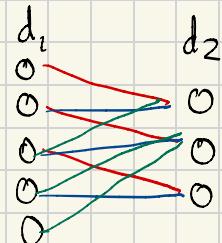
2. Shift invariance

## Convolutions (1D)



## Neural Net View

"kernel"  $w_1, w_2, w_3$

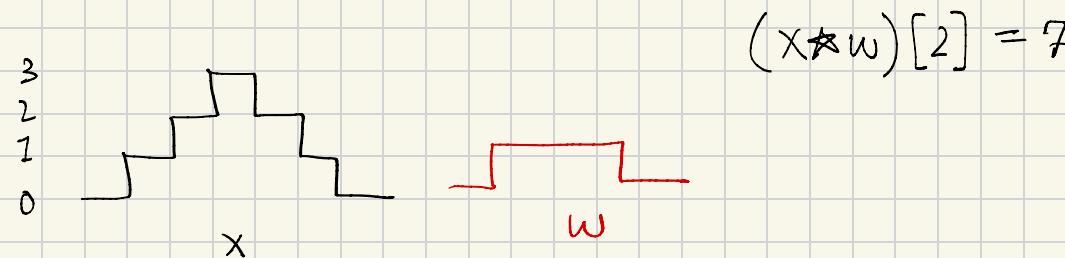


Weights are reused!

# weights = size(kernel)

vs

# weights =  $d_1 d_2$



## Convolutions (2D)

$$(x \star w)[2,2] = 2$$

$$x = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$w = \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix}$$

$$x \star w = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 2 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

$$(x \star w)[s,t] = \sum_{\sigma} \sum_{\tau} x[s+\sigma, t+\tau] w[\sigma, \tau]$$

$$x = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

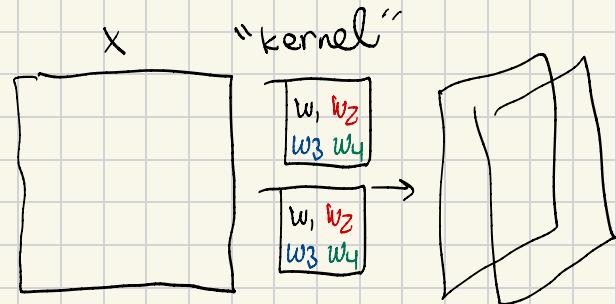
$$w = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

$$(x \star w)[2,2] = 3$$

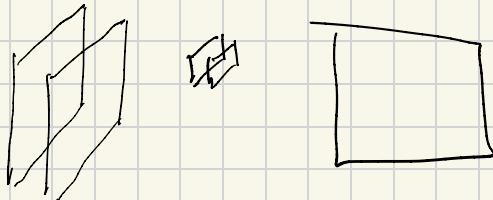
$$x \star w = \begin{pmatrix} & & \\ -3 & 3 & \\ & -3 & 3 \end{pmatrix}$$

## Neural Net View

$\text{Conv}(x, \text{kernel})$



## Convolutions (3D)



## More Considerations (ezyang)

- ↳ Padding
- ↳ Strides
- ↳ Pooling

## CNN Training

- ↳ keep track of reused weights

## Timeline of "Deep" Networks

Convolution → Residual → Transformer → Mixture of Experts

Deeper! Bigger! More data!

## Residual Networks (loss landscape residual)

Skipped connections

$$h^{(l)} = \sigma(W^{(l)} h^{(l-1)})$$

vs

$$h^{(l)} = \sigma(W^{(l)} h^{(l-1)}) + h^{(l-1)}$$

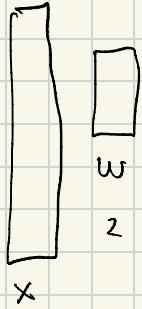
Pytorch Example (:

## Reminders

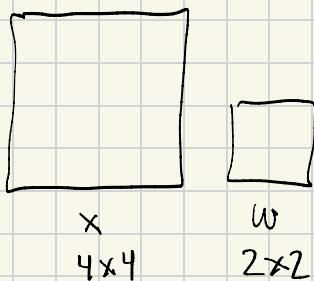
- ↳ Pset 7 due Wednesday
- ↳ Office hours
  - 10-11:30 Thursday 10/16
  - NO office hours Monday 10/20
- ↳ Midterm 10/21
- ↳ Practice exam!

## Review

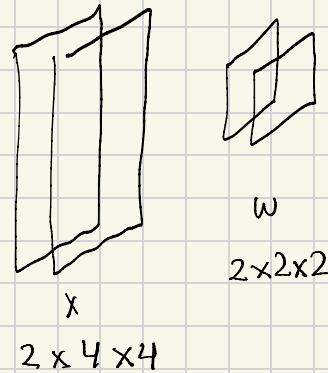
1D



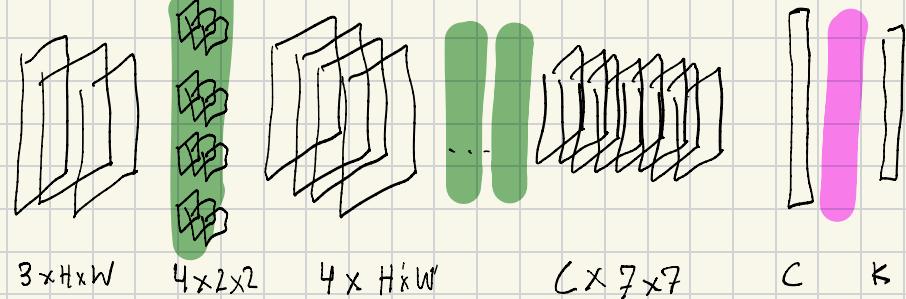
2D



3D



CNN



convolution  
Linear

## Motivation

Tabular data  
(e.g., IRIS, diabetes)

=> Linear

Image and audio data  
(e.g., MNIST, imagenet)

=> Convolution

Sequential data  
(e.g., text, time series)

=> Transformer

## Example

I live in Claremont. As usual, I will dress for —

BOS

TOKENIZATION

EOS



$x_1, x_2, \dots$

$x_n \in \mathbb{R}^d$

TRANSFORMER



$y_1, y_2, \dots$

$y_n \in \mathbb{R}^d$

## Transformer      Wishlist

1. handle variable length input

2. numeric representation of words

3. process context (e.g., I  $\Leftrightarrow$  Claremont  $\Leftrightarrow$  dress)

## Self Attention

Input:  $x_1, \dots, x_n$

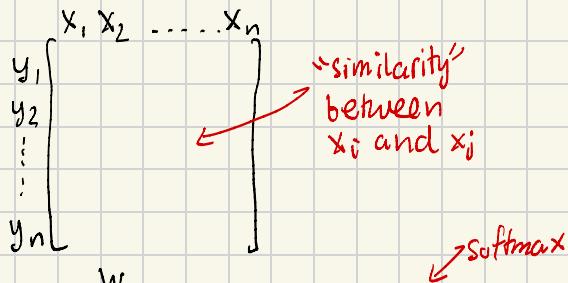
Output:  $y_1, \dots, y_n$

Idea: Build  $y_i$  out of  $x_1, \dots, x_n$

Attempt #1

group similar words

$$y_i = \sum_{j=1}^n w_{ij} x_j$$



$$w_{ij} = \frac{e^{<x_i, x_j>}}{\sum_k e^{<x_i, x_k>}}$$

But  $x_i$  appears three different places

1. query

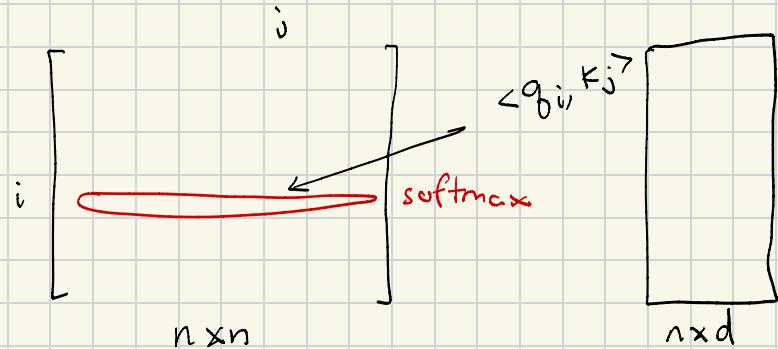
$$q_i = W^{(q)} x_i$$

2. key

$$k_i = W^{(k)} x_i$$

3. value

$$v_i = W^{(v)} x_i$$



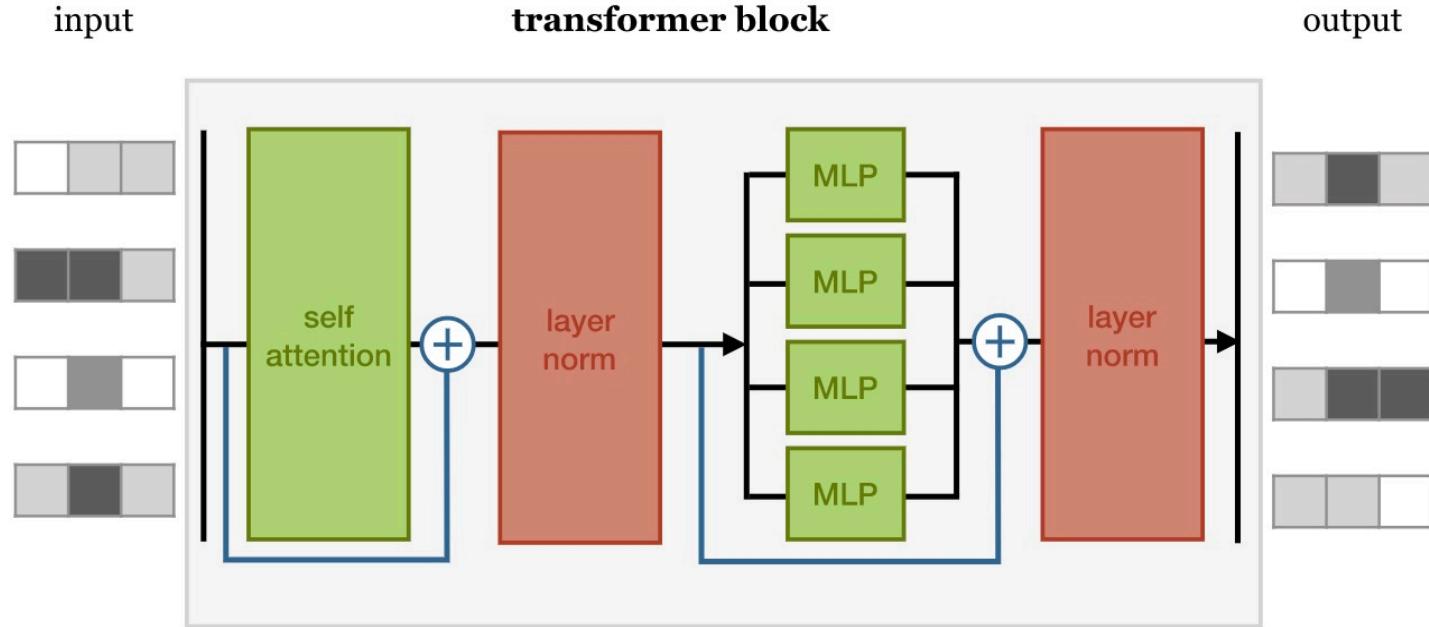
$$x \quad W^{(q)} \quad W^{(k)} \quad X^T$$

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

$$X \quad W^{(v)}$$

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

# Transformer



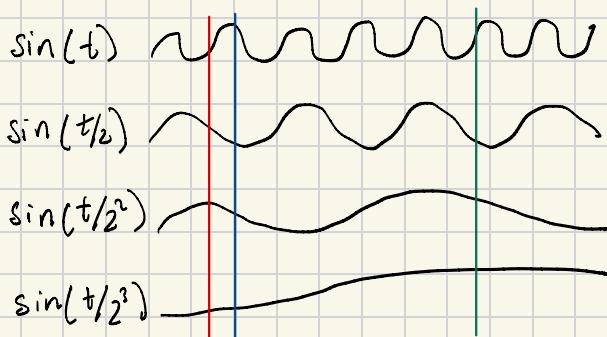
## Positional Encoding

"Jack gave water to Jill"

vs

"Jill gave water to Jack"

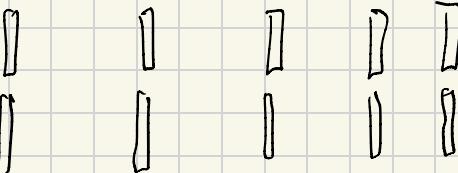
Idea: Add positional encoding to tokens



Jack gave water to Jill

Semantic  
embedding

Positional  
encoding



## Autoregressive Caching

	$x_1$	$x_2$	$\dots$	$x_n$	
$y_1$	0	0	0	0	
$y_2$	0	0	0		
$\vdots$	0	0			
$y_n$	0				

"causal mask"

$y_i$  built only from prior tokens

Benefit: Adding  $x_n$  ONLY changes  $y_n$

$\hookrightarrow O(n)$  updates vs  $O(n^2)$