

Week 7

10/7/2025

↳ Quiz

↳ Midterm 10/21

- practice exam(s) coming soon
- Linear Algebra → Transformers

Context:

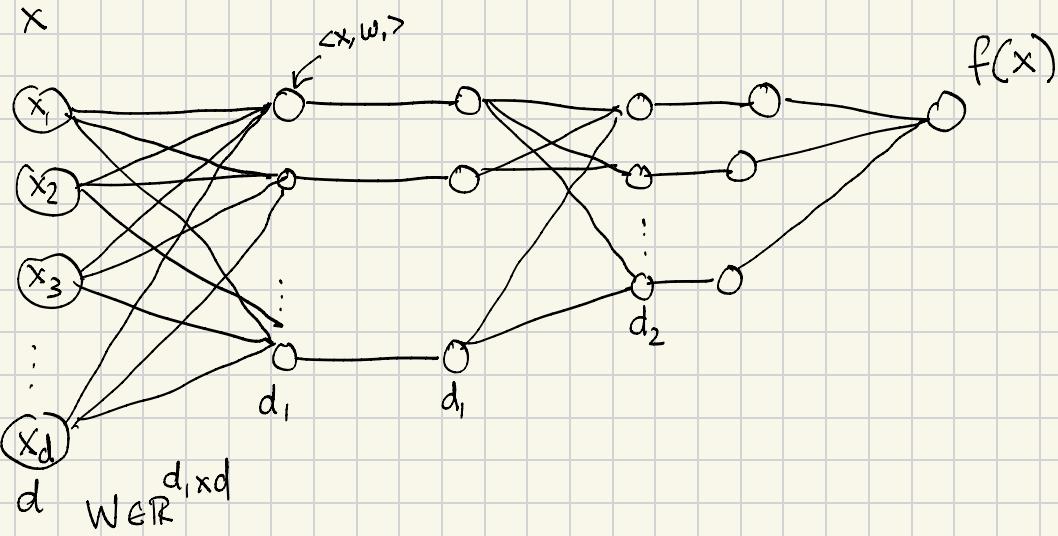
- Models {
- Linear Regression
 - Logistic Regression
 - Support Vector Machines

- Feature Selection {
- Feature Transformation
 - Kernel Trick
 - Reparameterization Trick

=> Neural Networks!

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

Review : Fully Connected NNs



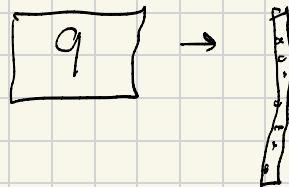
`nn.Linear(d, d1) nn.ReLU() ...`

$W_x = z$ $\sigma(z)$...

Issues with Linear Layers

1. Computationally expensive

2. Loss of "context"

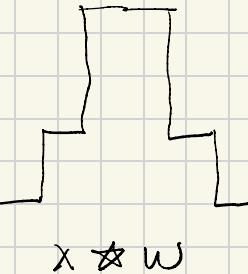
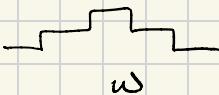
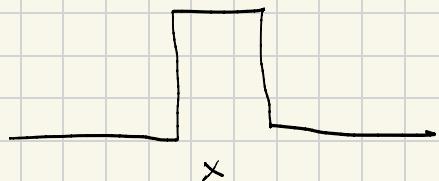


Wishlist for Convolutional Layers

1. Locality

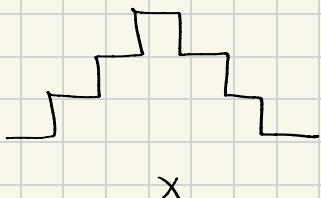
2. Shift invariance

Convolutions (1D)



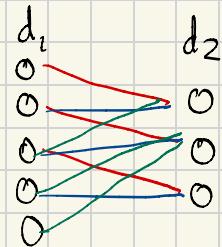
$$(x \star w)[t] = \sum_{\tau} x[t + \tau] w[\tau]$$

sanity ↴



Neural Net View

"kernel" w_1, w_2, w_3



Weights are reused!

weights = size(kernel)

vs

weights = $d_1 d_2$

Convolutions (2D)

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

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

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

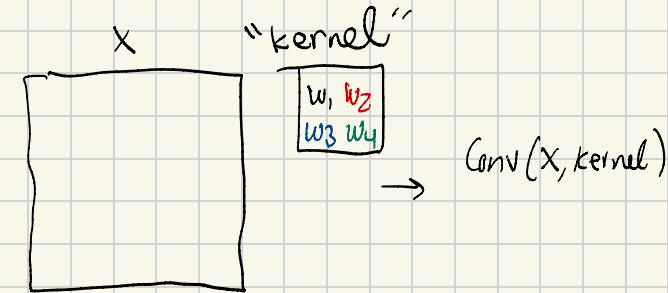
$$(x * w)[s, t] = \sum_{\sigma} \sum_{\tau} x[s + \sigma, t + \tau] w[\sigma, \tau]$$

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

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

$$x * w$$

Neural Net View



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 (: