

1. Science (Backpropagation, Regularization, SGD, GD, Vanishing and Exploding gradient problem, Little bit more on regularization)
2. Alchemy-like knowledge
(hidden units, how many layers, what kind of layers, why? What kind of architecture should we use for what kind of problem?
Not having a general theory, but figuring out by trial and error in the community.
)

Many kinds of Layers

↳ Linear Layer + Non linear activation functions

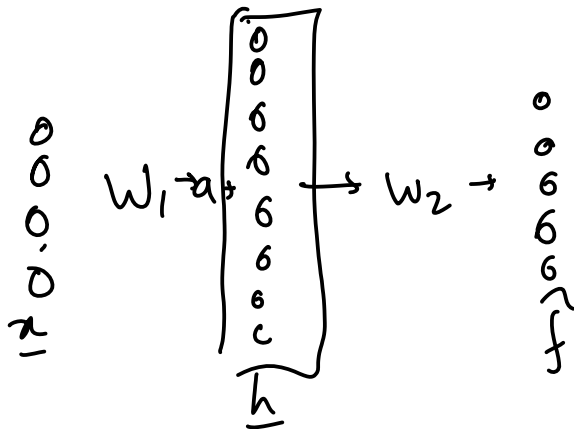
MLP: Multi-Layer Perceptron

Universal Approximation Theorem

2-layer MLP

$$\hat{f}(x) = W_2 \left(a \left(W_1 x + b_1 \right) \right) + b_2$$

number of hidden units

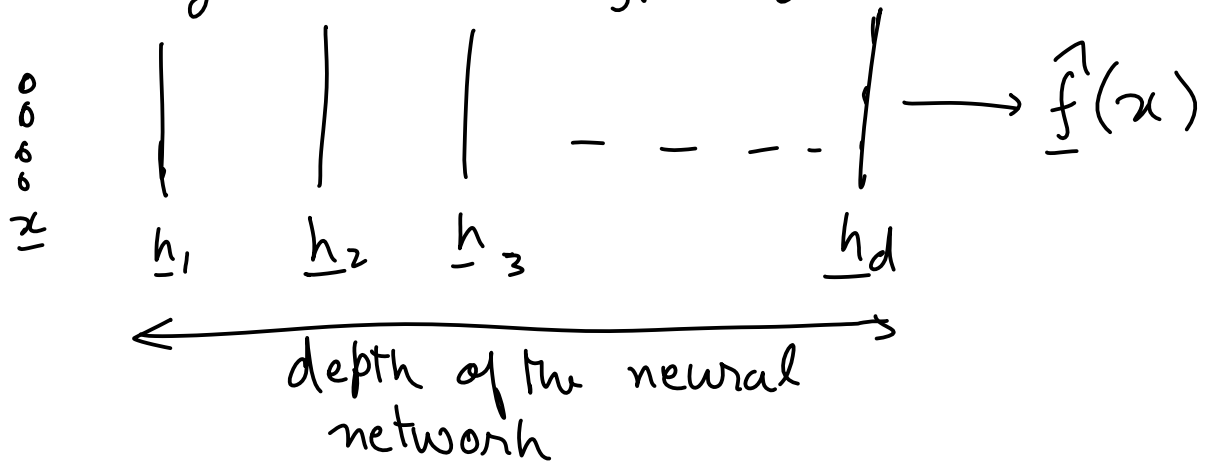


hidden ?
units ?

100, 1000, 10000, ∞

UAT: You can approximate any continuous function using an infinitely-wide 2-layer MLP with activation functions like, ReLU, sigmoid, tanh, ----

UAT': ∞ -layer MLP with sufficiently wide hidden unit

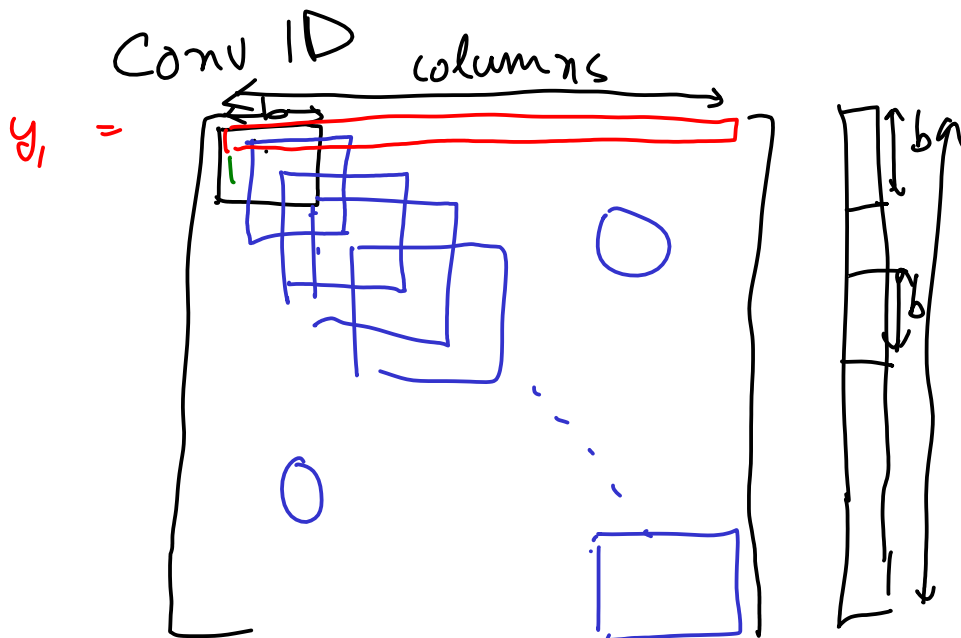


Fourier series } with ∞ expansion, you can approximate
Taylor series } any continuous function

① MLP layers
② Convolutional architecture, ③ Transformer architecture

↓
Mostly applied to images, data where there is repetition of patterns

Conv 1D is a particular kind of Linear Layer
Conv 2D



$$y_i = w_i \otimes x$$

$$y_i = \sum_{j \in I} w_{ij} x_{i-j}$$

W

learned

kernel

Conv 2D

← stride 200

200

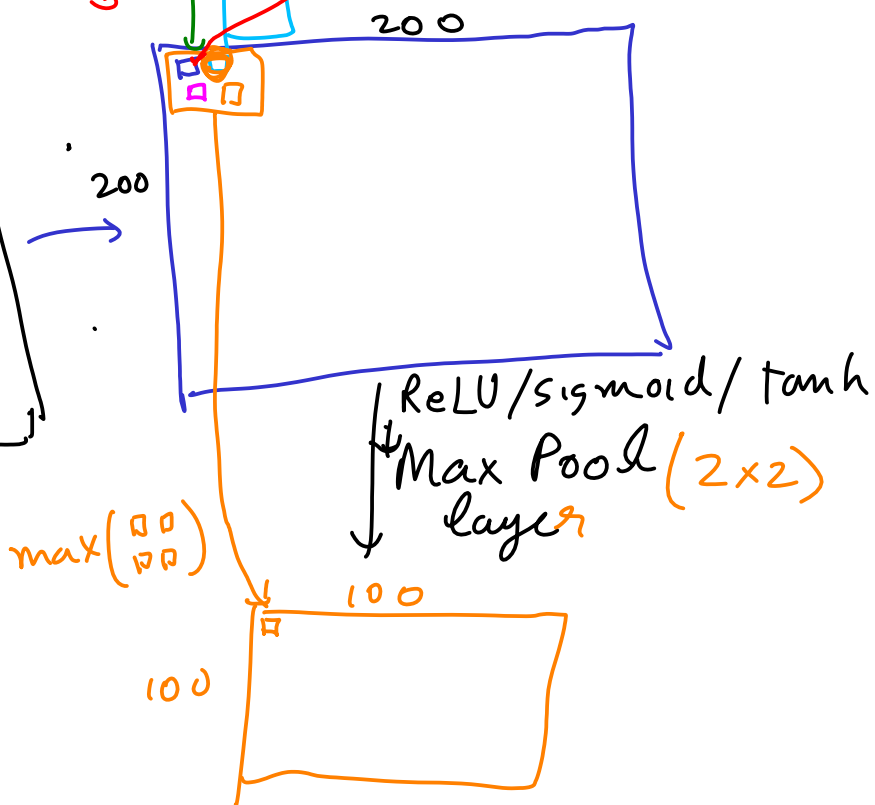
$$\begin{bmatrix} +1 & 0 & +1 \\ +1 & 0 & +1 \\ +1 & 0 & +1 \end{bmatrix}$$

edges in the image

$$\begin{bmatrix} +1 & 0 & +1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Sobel edge operator

$$\sum_{25} 5 \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} \times \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}$$



feature extractor
Pooling

④ Pooling → Max Pool 2D
→ Avg Pool 2D

③ Transformer / Self-attention / cross-attention

Hypernetwork

$$\underline{y} = a \left(\overset{\text{learnable}}{W_1} \underline{x} \right) \rightarrow \text{Linear layer}$$

output of another network?

learnable W_q, W_k, W_v

$$\underline{q} \underline{k}^T = \underbrace{W_q \underline{x} \underline{x}^T W_k^T}_{\text{O-1}}$$

softmax

$$\underline{q} = W_q \underline{x}$$

$$\underline{k} = W_k \underline{x}$$

$$\underline{v} = W_v \underline{x}$$

(self attention)

$$\text{softmax}(\underline{x})_i = \frac{\exp(x_i)}{\sum \exp(x)}$$

$$\begin{aligned} \underline{y} &= \text{softmax}(\underline{q} \underline{k}^T) \underline{v} \\ &= \underbrace{\text{softmax}(\underline{q} \underline{k}^T)}_W W_v \underline{x} \end{aligned}$$

Self-attention layer

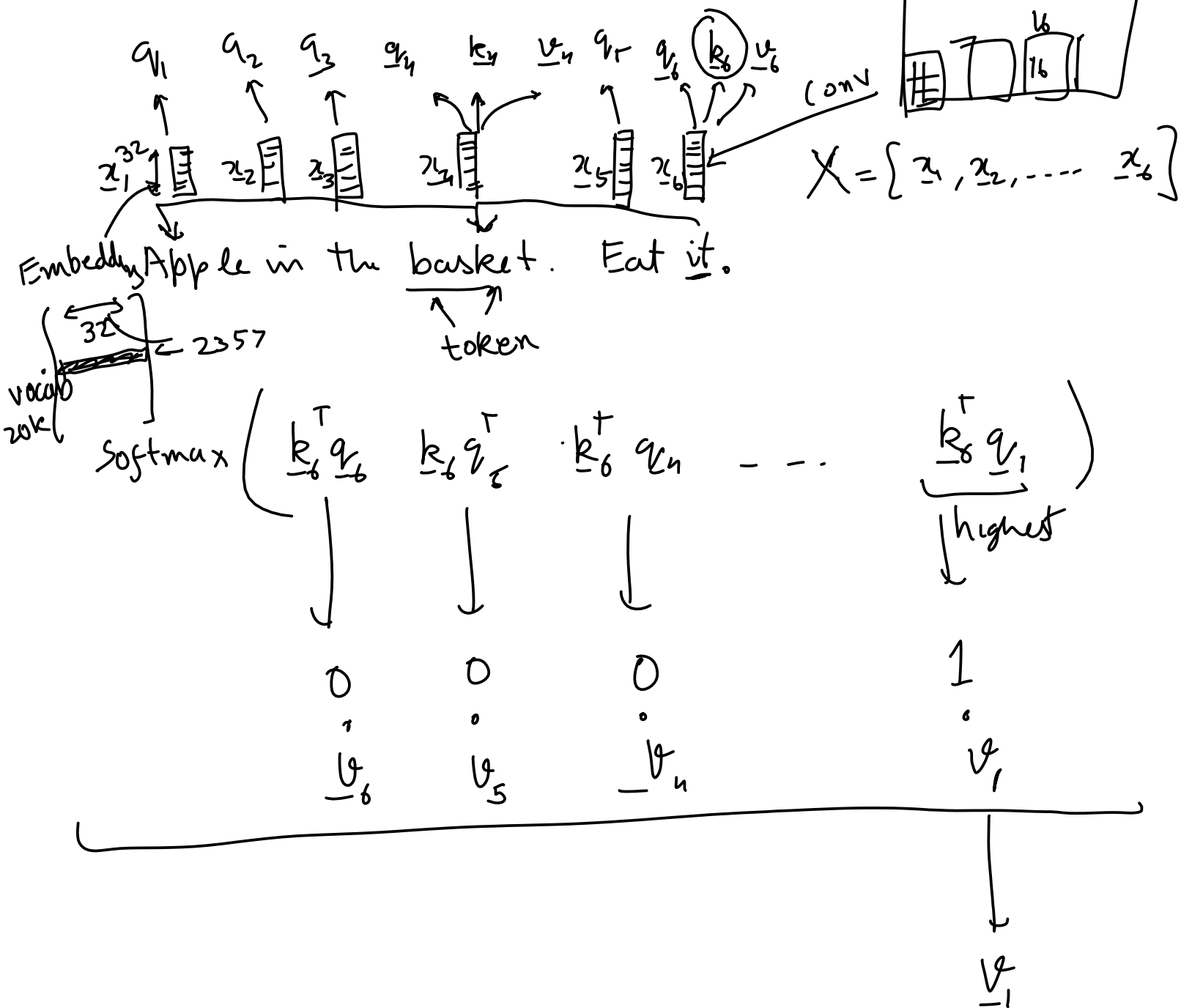
$$\underline{y} = \text{softmax}(\underline{Q} \underline{K}^T) \underbrace{W_v \underline{X}}_V$$

$$\underline{Q} = W_q \underline{X}$$

$$\underline{K} = W_k \underline{X}$$

$$\underline{V} = W_v \underline{X}$$

Natural language processing



Self attn₁

$X \rightarrow$ Self attn₂ } Multi head attn

① MLP

② Conv

③ Pooling

④ Transformer / self attn

Cross attn layer

when q, k come from diff data

than v

(5) Embedding Layer (indexing into a big matrix)