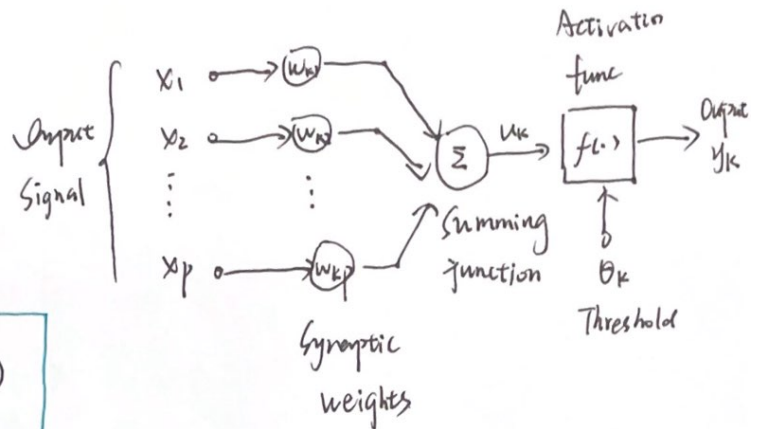
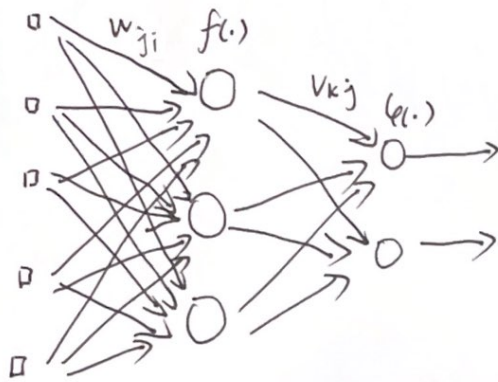


EE7403 LEC 14 Neural Networks and Deep Learning: From MLP to CNN

1. Neural Network and Deep CNN

① Network architecture of artificial neural networks (ANN), or simply, neural network.

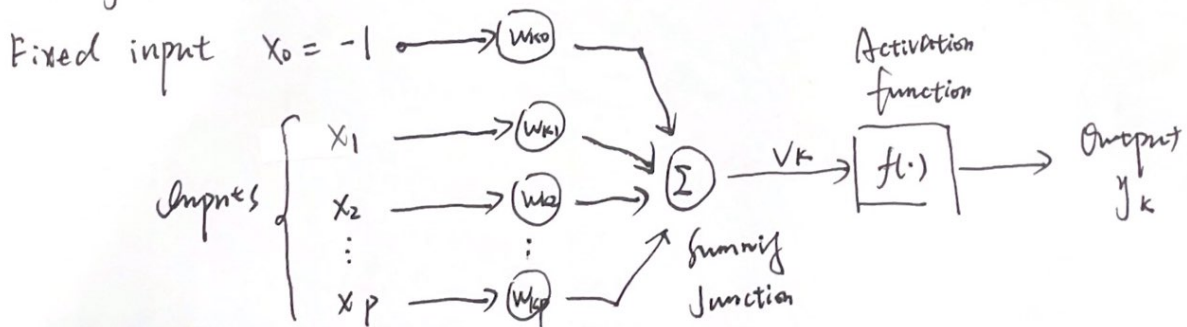


$$u_k = \sum_{j=1}^p w_{kj} x_j \quad y_k = f(u_k - \theta_k)$$

θ_k is an external parameter, we can consider

$$x_0 = 1, \quad w_{k0} = -\theta_k$$

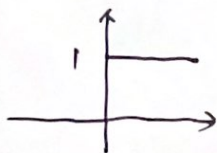
$$u_k = \sum_{j=0}^p w_{kj} x_j \quad y_k = f(u_k)$$



② Types of Activation Functions

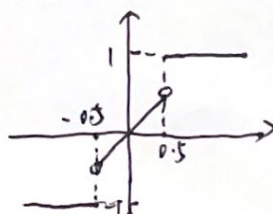
a. Binary function

$$f(u) = \begin{cases} 1, & u \geq 0 \\ 0, & u < 0 \end{cases}$$



b. Piece-wise-Linear

$$f(u) = \begin{cases} 1 & u \geq 0.5 \\ u & -0.5 \leq u < 0.5 \\ 0 & u \leq -0.5 \end{cases}$$



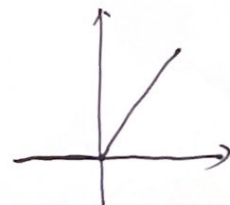
c. Sigmoid

$$f(u) = \frac{1}{1 + \exp(-au)}$$



d. ReLu

$$f(u) = \max(0, u) = \begin{cases} u, & u \geq 0 \\ 0, & u < 0 \end{cases}$$



2. ANN @ Multilayer Perceptron (MLP)

$$W_j = \begin{bmatrix} w_{j1} \\ w_{j2} \\ \vdots \\ w_{jn} \end{bmatrix} \quad V_k = \begin{bmatrix} v_{k1} \\ v_{k2} \\ \vdots \\ v_{kd} \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_d \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_c \end{bmatrix}$$

$$W = [w_1 \ w_2 \ \dots \ w_d] \quad \Rightarrow \quad Z_j = W_j^T X \quad \Rightarrow \quad Z = W^T X$$

$$V = [v_1 \ v_2 \ \dots \ v_c] \quad \text{similarly} \quad Y = V^T Z$$

$$Z_j = f\left(\sum_{i=1}^n w_{ji} x_i\right) = f(W_j^T X)$$

$$Y_k = f\left(\sum_{j=1}^d v_{kj} Z_j\right) = f(V_k^T Z)$$

$$\Downarrow = f\left[\sum_{j=1}^d v_{kj} f\left(\sum_{i=1}^n w_{ji} x_i\right)\right]$$

$$Y = f(V^T Z) = f[V^T f(W^T X)]$$

If f is a linear function: $Y = f[V^T f(W^T X)] = V^T W^T X = (WV)^T X = U^T X$
 just same as single layer

where $U = WV$

To design a multi-layer network, the activation function should be nonlinear ★

② Single-Layer Neural Network

$$Y = W^T X \quad \text{target output: } t$$

error: $e = t - y = t - W^T X$

The mean square error: $J(W) = E\{e^2\} = E\{(t - y)^2\} = E\{(t - W^T X)^2\}$ $W = \frac{t}{X}$

least square method: quadratic:

$$\frac{\partial E\{e^2\}}{\partial W} = \frac{\partial (t^2 - 2W^T E\{xt\} + W^T E\{xx^T\} W)}{\partial W} = 0 \Rightarrow E\{xx^T\} W = E\{xt\}$$

$$E\{xx^T\} \approx \frac{1}{q} \sum_{i=1}^q x_i x_i^T \quad E\{xt\} \approx \frac{1}{q} \sum_{i=1}^q x_i t_i$$

Let $X = [x_1 \ x_2 \ \dots \ x_q]$ $t = [t_1 \ t_2 \ \dots \ t_q]$

$$q \cdot E\{xx^T\} \approx \sum_{i=1}^q x_i x_i^T = X X^T \quad q \cdot E\{xt\} \approx \sum_{i=1}^q x_i t_i = X t^T$$

$$E\{xx^T\} W = E\{xt\} \Rightarrow X X^T W = X t^T \Rightarrow W = (X X^T)^{-1} X t^T$$

EE7403 LEC14

gradient descent method:

$$W \leftarrow W + \Delta W$$

$$\Delta W = -\eta \nabla J(W)$$

$$\nabla J(W) = \frac{\partial J(W)}{\partial W} = \frac{\partial E\{(t - W^T x)^2\}}{\partial W} = -2E\{(t - W^T x)x\}$$

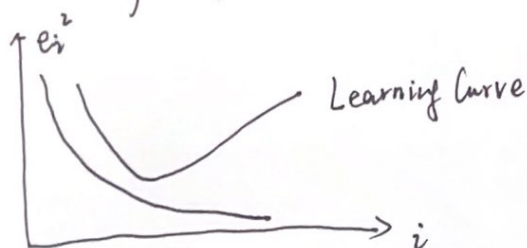
$$= -2E\{e x\} \cong -2e_i x_i$$

$$W \leftarrow W + \eta e_i x_i = W + \eta (t_i - y_i) x_i$$

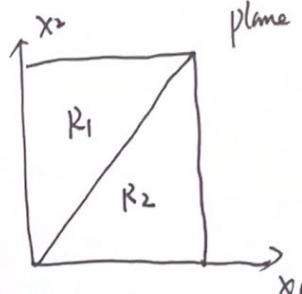
$$\rightarrow W = (X X^T)^{-1} X t^T$$

③ Learning Curve and Decision Boundary

$$W \leftarrow W + \eta (t_i - y_i) x_i$$



decision boundary is a straight line or plane or hyper-plane



④ backpropagation

$$y = f[V^T f(W^T x)]$$

$$J(W, V) = E\{e^2\} = E\{\|t - y\|^2\} = E\{\|t - f[V^T f(W^T x)]\|^2\} \quad (q_k = \sum_{j=1}^d v_{kj} z_j)$$

$$= \frac{1}{2} \sum_{k=1}^c [t_k - f(q_k)]^2 = \frac{1}{2} \sum_{k=1}^c [t_k - f(\sum_{j=1}^d v_{kj} z_j)]^2$$

$$W \leftarrow W - \eta \nabla J(W)$$

$$\frac{\partial J(W, V)}{\partial v_{kj}} = \frac{\partial J}{\partial q_k} \cdot \frac{\partial q_k}{\partial v_{kj}} = -(t_k - y_k) \cdot f'(q_k) z_j$$

$$V \leftarrow V - \eta \nabla J(V)$$

Similarly, $J(W, V) = E\{\|t - y\|^2\} = \frac{1}{2} \sum_{k=1}^c [t_k - f(q_k)]^2 = \frac{1}{2} \sum_{k=1}^c [t_k - f(\sum_{j=1}^d v_{kj} z_j)]^2$

$$= \frac{1}{2} \sum_{k=1}^c [t_k - f(\sum_{j=1}^d v_{kj} f(q_j))]^2 = \frac{1}{2} \sum_{k=1}^c [t_k - f(\sum_{j=1}^d v_{kj} f(\sum_{i=1}^n w_{ji} x_i))]^2$$

$$\frac{\partial J(W, V)}{\partial w_{ji}} = \frac{\partial J}{\partial z_j} \cdot \frac{\partial z_j}{\partial q_j} \cdot \frac{\partial q_j}{\partial w_{ji}} = -\left[\sum_{k=1}^c (t_k - y_k) \cdot f'(q_k) v_{kj}\right] f'(q_j) x_i$$

$$v_{kj} \leftarrow v_{kj} - \eta \frac{\partial J(W, V)}{\partial v_{kj}} \quad w_{ji} \leftarrow w_{ji} - \eta \frac{\partial J(W, V)}{\partial w_{ji}}$$

$$v_{kj} \leftarrow v_{kj} + \eta (t_k - y_k) \cdot f'(q_k) z_j$$

$$w_{ji} \leftarrow w_{ji} + \eta \left[\sum_{k=1}^c (t_k - y_k) f'(q_k) v_{kj}\right] f'(q_j) x_i$$

$$f(q) = \frac{1}{1 + \exp(-aq)} = \frac{\exp(aq)}{1 + \exp(aq)} = a f(q) [1 - f(q)]$$

propagation
procedure:

Assume n training samples, no prior info.

(1) Initialization

(2) Presentation of training samples should be randomized from epoch to epoch

(3) Forward computation

(4) Backward computation

(5) Iteration. criterion can be the number of iterations or the rate of change of the average error small enough

Stop training at minimum of the error on the validation set.

x overfit

3. Problems of Machine Learning

Problems of local minima (under-fitting). poor generalization (overfitting)

① Why deep layer though 2 layers are enough:

- deep layer NN is easier to train. a shallow NN with many neurons is hard to optimize. fewer neurons per layer are often with better gradient flow.
- deep layer needs fewer parameters, there're fewer neuron per layer reuse parameters, making it more efficient.
- learns hierarchical features Better generalization. deep layer structure like feature extractor.

② $\min_f \hat{e}^2 = \min_{\hat{f}} \sum_i \|y_i - \hat{f}(x_i)\|_2^2 \Rightarrow 0$, can only find $y_i = \hat{f}(x_i)$, not $y = f(x)$.

how to make $\hat{f}(x) \Rightarrow f(x)$? Regulation.

4. Convolutional Neural Networks, CNN

① different from MLP.

1) network architecture

2) Convolution

3) simple nonlinear activation function ReLU.

4) pooling

5) deep. large number of layer.

Input size $P \times Q$. C channels.Output size $P \times Q$. D channels.

$$x_{i,j,k}, \quad 1 \leq i \leq P, \quad 1 \leq j \leq Q, \quad 1 \leq k \leq C.$$

$$y_{i,j,k}, \quad 1 \leq i \leq P, \quad 1 \leq j \leq Q, \quad 1 \leq k \leq D$$

$$y_{i,j,k} = \sum_{n=1}^C \sum_{m=1}^P \sum_{l=1}^Q w_{l,m,n,k} \cdot x_{i-l,j-m,n} + b_k \quad 1 \leq i \leq P, \quad 1 \leq j \leq Q, \quad 1 \leq k \leq D$$

input channel

CNN capture the local feature. includes a bias term

Along channel dim, all inputs of different channels are all fully connected.

② Characteristics

1) Local Perceptive Fields

2) Shared Weights (gives translation invariance. num of parameters ↓)

3) Full channel connected.

4) Bias Term

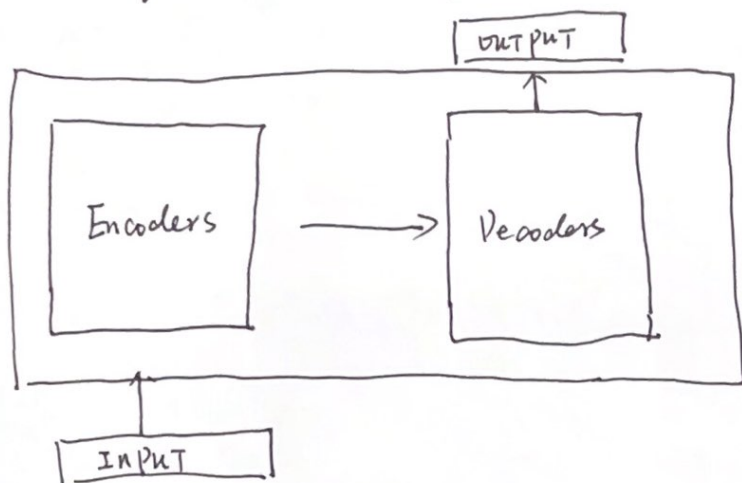
5) Hierarchy Features

A CNN is essentially a regularized version of MLP.

EE7403 LEC 15 Deep Learning: From CNN to Transformer.

Attention is all you need.

Transformer captures the relationships between each word/token in a sequence with every other word/token.



The Transformer.

Encoders + Decoders

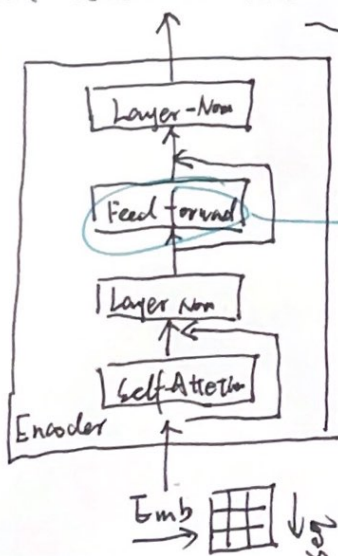
1. Embedding

Embed each input token into a feature vector.

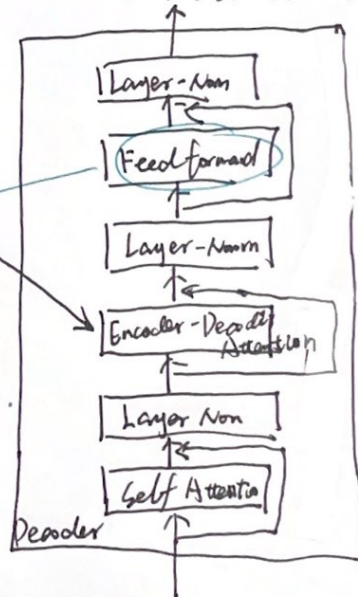
1 token is a vector

2. Encoder & Decoder

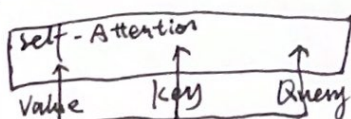
Both Attention and Feed-forward layers have a residual skip-connection



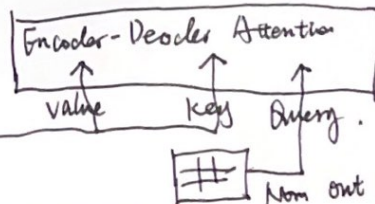
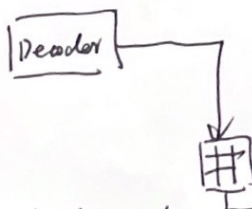
a couple of linear layer with a ReLU activation in between.



3. Attention



(Input embeds and position encodings)



Three copies of each word/token are generated for self-attention by linear propagation.

trainable W_q, W_k, W_v

$$\text{Let } X = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}, Q = \begin{pmatrix} Q_1 \\ Q_2 \\ Q_3 \\ Q_4 \end{pmatrix}, K = \begin{pmatrix} K_1 \\ K_2 \\ K_3 \\ K_4 \end{pmatrix}, V = \begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{pmatrix}$$

Then: $Q = XW_q, K = XW_k, V = XW_v$ each 'row' of these matrices corresponds

① $R = QK^T$ the relevance between words. to one token in the source sequence

Dot product generates similarity between words.

used as factors

$$\textcircled{2} Z = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad Q \cdot K \cdot V \cdot Z \in \mathbb{R}^{n \times k}$$

dimension k

$$\text{Attention} = W = QK^T$$

Output $Z = WV$ weighted sum of all input tokens. Capture global info.

\swarrow
x learnt from training data but generated by specific test input.

Attention itself is learnable.

4. Learnable feed forward networks

$$Y = h(h(h(XW_1)W_2)W_3)$$

$$(n \times k)(k \times k) \Rightarrow (n \times k)$$

all tokens in inputs/outputs share same learnable parameters
each single learnable parameter go through all points