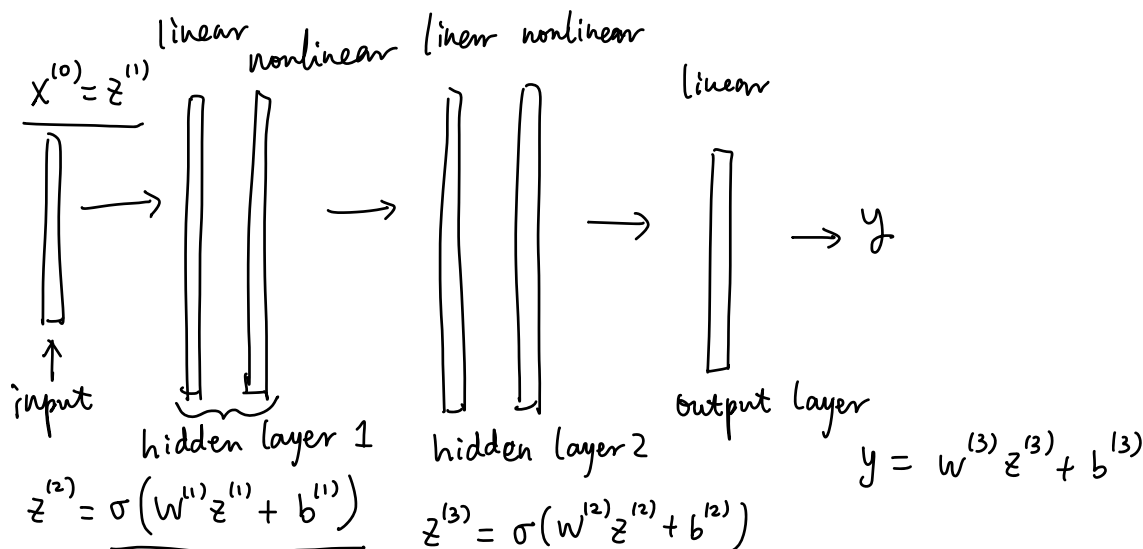


Review:

Feedforward neural network (FNN):



example: take FNN with 1 hidden layer. activation function

$$W^{(1)} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \quad b^{(1)} = \begin{pmatrix} -2 \\ 0 \end{pmatrix}$$

$$\text{output layer: } W^{(2)} = (1, 1), \quad b^{(2)} = 0$$

$$\text{input: } X^{(0)} = z^{(1)}$$

$$\text{hidden layer: } z^{(2)} = \sigma(W^{(1)}z^{(1)} + b^{(1)})$$

$$= \sigma\left(\begin{pmatrix} 2 \\ 1 \end{pmatrix} z^{(1)} + \begin{pmatrix} -2 \\ 0 \end{pmatrix}\right)$$

$$= \sigma\left(\begin{pmatrix} 2z^{(1)} - 2 \\ z^{(1)} \end{pmatrix}\right)$$

$$= \text{ReLU}\left(\begin{pmatrix} 2z^{(1)} - 2 \\ z^{(1)} \end{pmatrix}\right) = \begin{pmatrix} \text{ReLU}(2z^{(1)} - 2) \\ \text{ReLU}(z^{(1)}) \end{pmatrix}$$

$$\text{output layer: } y = W^{(2)}z^{(2)} + b^{(2)}$$

$$\sigma(x) = \text{ReLU}(x)$$

$$= \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

$$= \max\{x, 0\}$$

$$\begin{aligned}
&= (1, 1) \bar{z}^{(2)} + 0 \\
&= (1, 1) \begin{pmatrix} \text{ReLU}(2\bar{z}^{(1)} - 2) \\ \text{ReLU}(\bar{z}^{(1)}) \end{pmatrix} + 0 \\
&= \text{ReLU}(2\bar{z}^{(1)} - 2) + \text{ReLU}(\bar{z}^{(1)})
\end{aligned}$$

FNN: $y = \text{ReLU}(2x^{(0)} - 2) + \text{ReLU}(x^{(0)})$

$$y = \max\{2x^{(0)} - 2, 0\} + \max\{x^{(0)}, 0\}$$

(1) if $x < 0$.

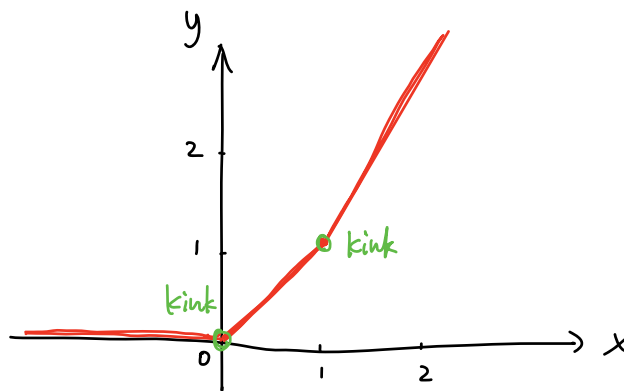
$$2x - 2 < 0, \quad x < 0 \quad \Rightarrow \quad y = 0 + 0 = 0$$

(2) if $0 < x < 1$.

$$2x - 2 < 0, \quad x > 0 \quad \Rightarrow \quad y = 0 + x = x$$

(3) if $x > 1$.

$$2x - 2 > 0, \quad x > 0 \quad \Rightarrow \quad y = 2x - 2 + x = 3x - 2$$



Remark: (1) FNN is a function

(2) The function in this example has two kinks

(3) FNN with 1 hidden layer with 2 neurons : two kinks

FNN with 2 hidden layer with 2 neurons : four kinks

example: FNN with 2 hidden layer with 2 neurons each layer (HW)

$$\text{input: } x^{(0)} = z^{(1)}$$

$$\text{hidden layer 1: } z^{(2)} = \sigma(w^{(1)} z^{(1)} + b^{(1)})$$

$$\text{hidden layer 2: } z^{(3)} = \sigma(w^{(2)} z^{(2)} + b^{(2)})$$

$$\text{output: } y = w^{(3)} z^{(3)} + b^{(3)}$$

expressivity of neural network is increasing when we use more hidden layers and more neurons

Definition (FNN): An FNN is a function $f: \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$. For an input $x \in \mathbb{R}^{d_{\text{in}}}$, the output $y \in \mathbb{R}^{d_{\text{out}}}$ is given by:

$$\begin{cases} z^{(1)} = x \\ z^{(i+1)} = \sigma(w^{(i)} z^{(i)} + b^{(i)}), & i = 1, 2, \dots, L \\ y = w^{(L+1)} z^{(L+1)} + b^{(L+1)} \end{cases}$$

Here, we say that the number of hidden layers is L . The number of neurons in each hidden layer is denoted by H_i for $1 \leq i \leq L$.

Therefore, the parameters have the dimension:

$$w^{(1)} \in \mathbb{R}^{H_1 \times d_{\text{in}}}, \quad b^{(1)} \in \mathbb{R}^{H_1}$$

$$w^{(2)} \in \mathbb{R}^{H_2 \times H_1}, \quad b^{(2)} \in \mathbb{R}^{H_2}$$

\vdots

$$W^{(L)} \in \mathbb{R}^{H_L \times H_{L-1}} \quad b^{(L)} \in \mathbb{R}^{H_L}$$

$$W^{(L+1)} \in \mathbb{R}^{d_{out} \times H_L} \quad b^{(L+1)} \in \mathbb{R}^{d_{out}}$$

Example: An FNN with $d_{in} = d_{out} = 1$. 3 hidden layers.

5 neurons in each hidden layer ($H_1 = H_2 = H_3 = 5$)

total number of parameters in this FNN?

$$W^{(1)} \in \mathbb{R}^{5 \times 1}, \quad b^{(1)} \in \mathbb{R}^5 \quad \rightarrow \quad 5 + 5 = 10$$

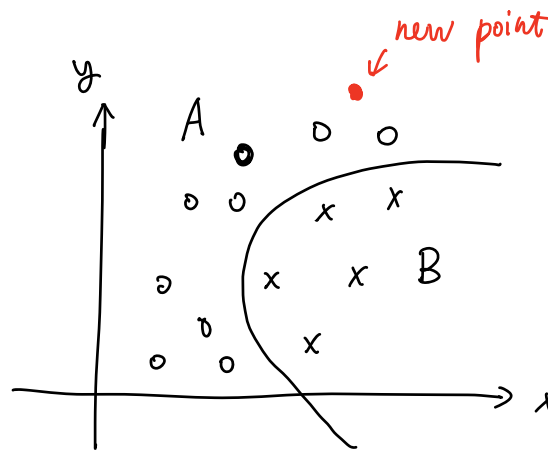
$$W^{(2)} \in \mathbb{R}^{5 \times 5}, \quad b^{(2)} \in \mathbb{R}^5 \quad \rightarrow \quad 25 + 5 = 30$$

$$W^{(3)} \in \mathbb{R}^{5 \times 5}, \quad b^{(3)} \in \mathbb{R}^5 \quad \rightarrow \quad 25 + 5 = 30$$

$$W^{(4)} \in \mathbb{R}^{1 \times 5}, \quad b^{(4)} \in \mathbb{R}^1 \quad \rightarrow \quad 5 + 1 = 6$$

$$\text{total number: } 10 + 30 + 30 + 6 = \underline{76}$$

classification problem:



$$\text{FNN: } y = \sigma(W^{(3)} \sigma(W^{(2)} \sigma(W^{(1)} x + b^{(1)}) + b^{(2)}) + b^{(3)})$$

Find out $(W^{(1)}, W^{(2)}, W^{(3)}, b^{(1)}, b^{(2)}, b^{(3)})$ such that

$$y(x^{(i)}) = \begin{cases} (1, 0) & \text{if } x^{(i)} \in A \\ (0, 1) & \text{if } x^{(i)} \in B \end{cases}$$

Define the loss function (or called cost function) :

$$L(w^{(1)}, w^{(2)}, w^{(3)}, b^{(1)}, b^{(2)}, b^{(3)}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|y^{(i)} - y(x^{(i)})\|_2^2$$

Training neural network (optimization problem) :

Find $w^{(1)}, w^{(2)}, w^{(3)}, b^{(1)}, b^{(2)}, b^{(3)}$, such that the loss function is minimized.

$f(x)$ take the min value at $x=x^*$.

↓

optimization : minimizing $f(x)$. $x^* = \operatorname{argmin} f(x)$

example: $f(x) = x^2$. take the min at $x=0$.

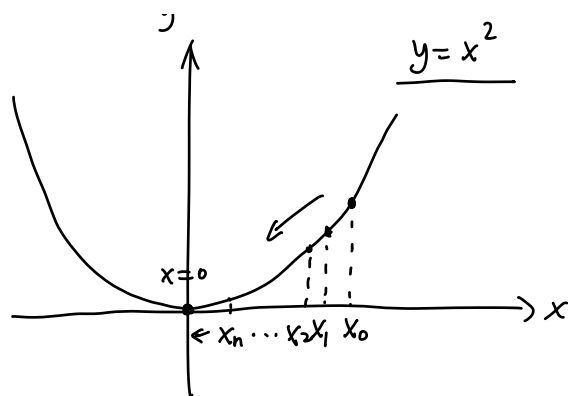
Gradient descent : find out min of $f(x)$

$$\begin{aligned} f(x+\Delta x) &= f(x) + \Delta x \cdot f'(x) + O(\Delta x^2) \leftarrow \text{Taylor expansion} \\ &\approx f(x) + \Delta x \cdot f'(x) \end{aligned}$$

take $\Delta x = -\eta \cdot f'(x)$. $\eta > 0$

$$\begin{aligned} f(x+\Delta x) &\approx f(x) - \eta \cdot f'(x) \cdot f'(x) \\ &= f(x) - \eta \cdot (f'(x))^2 \leq f(x) \end{aligned}$$

$$f(\underbrace{x - \eta \cdot f'(x)}) \leq f(x)$$



$$\begin{aligned} x_1 &= x_0 - \eta \cdot f'(x_0) \\ &= x_0 - \eta \cdot 2x_0 \end{aligned}$$

$$\underline{x_2 = x_1 - \eta \cdot f'(x_1)}$$

take the initial guess $x = x_0$, do the iteration

$$x_{n+1} = x_n - \underbrace{\eta}_{\uparrow} \cdot f'(x_n).$$

learning rate — small positive constant

generalization to vector case: vector vector

$$\begin{aligned} f(x + \Delta x) &= f(x) + \underbrace{\Delta x} \cdot \underbrace{\nabla f(x)} + O(\Delta x^2) \\ &\approx f(x) + \Delta x \cdot \nabla f(x) \end{aligned}$$

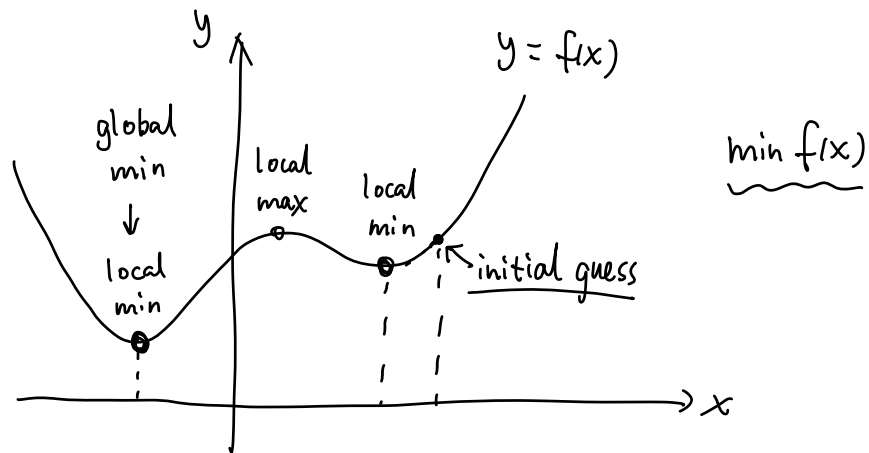
take $\Delta x = -\eta \nabla f(x)$. η small positive constant

$$\begin{aligned} f(x + \Delta x) &\approx f(x) + (-\eta \nabla f(x) \cdot \nabla f(x)) \\ &= f(x) - \eta |\nabla f(x)|^2 \leq f(x) \end{aligned}$$

$$f(x - \eta \nabla f(x)) \leq f(x)$$

take the initial guess $x = x_0$. do the iteration

$$x_{n+1} = x_n - \eta \nabla f(x_n)$$



critical point: $f'(x) = 0$

local min: critical point x_0 + $f''(x_0) > 0$

global min: min in the local min

gradient decent may not always find global minimum (HW)