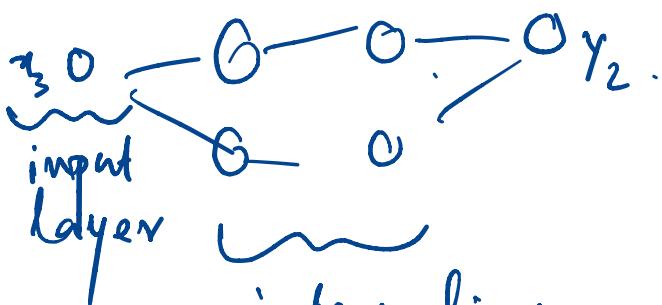
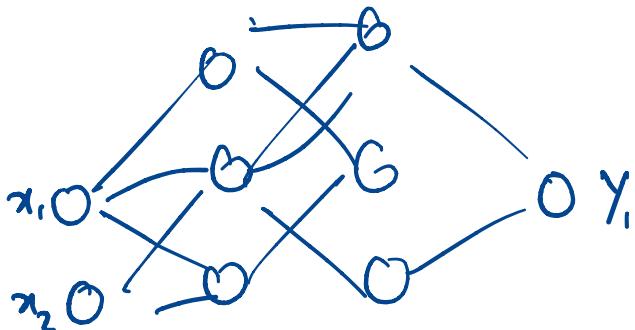


Lecture 3.



w
weighted
linear
combⁿ



$$y = f_{\theta}(x)$$

intermediate
layers.

weights & biases.

Data: $(x^1, y^1), \dots, (x^N, y^N)$

Find θ , that minimizes $\frac{1}{N} \sum_{i=1}^N \text{loss}(f_{\theta}(x^i), y^i) := E(\theta)$

minimize $E(\theta)$

θ :



Setting $\nabla E(\theta) = 0$

... solve θ^* .

Think line fitting.

Closed-form
expressⁿ!

Computational procedure

$-\nabla E(\theta)$ { ... gradient descent.

... denotes
the direction
of steepest
descent.

... $\theta^\circ \xrightarrow{\text{walk along}} -\nabla E(\theta^\circ)$

① how far
to walk?

walk
along $-\nabla E(\theta^\circ) \xrightarrow{\downarrow} \theta'$

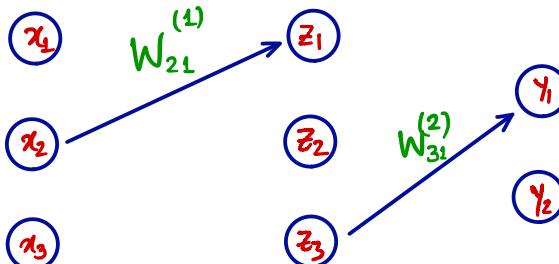
② when
to stop?

! step-size / learning
rate

• $\theta^{k+1} \leftarrow \theta^k - \eta \frac{\nabla E(\theta^k)}{\text{Compute.}}$

Non-Linear regression via neural networks (NN)

Neural networks provide a parametric map between input and output vectors in terms of weights and biases between each of its layers. Consider a NN which has $x \in \mathbb{R}^3$ as inputs and $y \in \mathbb{R}^2$ as outputs, and has one hidden layer with 3 neurons. Let's draw the NN architecture.



$$x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \in \mathbb{R}^3$$

$$z = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} \in \mathbb{R}^3$$

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \in \mathbb{R}^2$$

The functional relationships are given by

$$z = \varphi^{(1)} (W^{(1)} x + b^{(1)}),$$

$$y = \varphi^{(2)} (W^{(2)} z + b^{(2)}).$$

φ 's are activation functions.

$$z = \varphi^{(1)}(W^{(1)}x + b^{(1)}),$$

$$y = \varphi^{(2)}(W^{(2)}z + b^{(2)}).$$

$W^{(1)}$ and $W^{(2)}$ are weights.
 $b^{(1)}$ and $b^{(2)}$ are biases.

- Dimensions of W 's and b 's ?
 $W^{(1)} \in \mathbb{R}^{3 \times 3}$, $W^{(2)} \in \mathbb{R}^{2 \times 3}$,
 $b^{(1)} \in \mathbb{R}^3$, $b^{(2)} \in \mathbb{R}^2$.
- How do we tune the weights and biases ?

Suppose you are given a pair (x, t) , where $x \in \mathbb{R}^3$ and $t \in \mathbb{R}^2$. We aim to find $\theta := (W^{(1)}, W^{(2)}, b^{(1)}, b^{(2)})$ such that the error between t and y is minimized.

Consider an error (aka loss function) defined as

$$E(\theta) = \frac{1}{N} \sum_{i=1}^N E^{(i)}(\theta) \quad \parallel E^{(i)}(\theta) := \frac{1}{2} \left[(t_1^{(i)} - y_1)^2 + (t_2^{(i)} - y_2)^2 \right].$$

To minimize E , start with some θ and

take a step along the negative gradient of E , i.e., update θ as

$$\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} E(\theta)$$

for a suitable step-length η .

- How do we calculate $\nabla_{\theta} E(\theta)$?

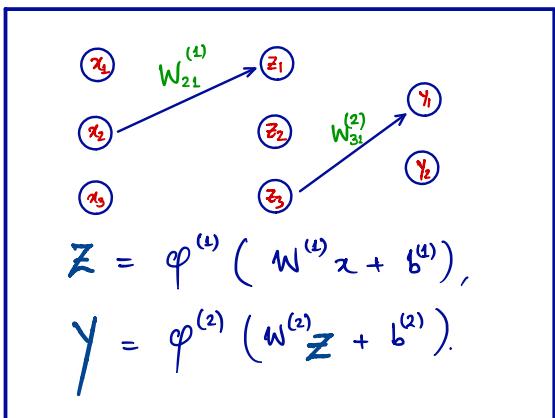
Unpacking the gradient, we seek to calculate

$$\frac{\partial E}{\partial W_{11}^{(1)}}, \frac{\partial E}{\partial W_{12}^{(1)}},$$

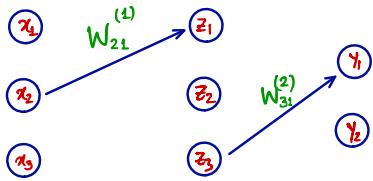
$$\dots \frac{\partial E}{\partial W_{11}^{(2)}}, \frac{\partial E}{\partial W_{12}^{(2)}},$$

$$\dots \frac{\partial E}{\partial b_1}, \dots$$

$$\frac{\partial E}{\partial b_1^{(2)}}, \dots \frac{\partial E}{\partial b_2^{(2)}}.$$



We calculate these derivatives beginning from the last layer first, and proceed backwards.



$$z = \varphi^{(1)}(W^{(1)}x + b^{(1)}),$$

$$y = \varphi^{(2)}(W^{(2)}z + b^{(2)}).$$

"Let's calculate $\frac{\partial E}{\partial W_{31}^{(2)}}$.

New notation:

$$z^{in} = W^{(1)}x + b^{(1)}$$

$$y^{in} = W^{(2)}z + b^{(2)}.$$

$$\therefore z = \varphi^{(1)}(z^{in}) \text{ and } y = \varphi^{(2)}(y^{in})$$

Then, we have

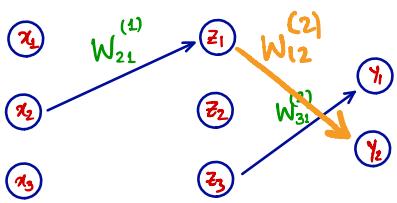
$$\frac{\partial E}{\partial W_{31}^{(2)}} = \underbrace{\frac{\partial E}{\partial y_1}}_{\text{Call it } \delta_1^{(2)}} \cdot \underbrace{\frac{\partial y_1}{\partial y_1^{in}}}_{=z_3} \cdot \underbrace{\frac{\partial y_1^{in}}{\partial W_{31}^{(2)}}}_{= \delta_1 \cdot z_3}$$

Call it $\delta_1^{(2)}$

$$y_1^{in} = z_1 \cdot W_{11}^{(2)} + z_2 \cdot W_{21}^{(2)}$$

Also, we have

$$\frac{\partial E}{\partial b_L^{(2)}} = \underbrace{\frac{\partial E}{\partial y_1}}_{= \delta_1^{(2)}} \cdot \underbrace{\frac{\partial y_1}{\partial y_1^{in}}}_{= 1} \cdot \underbrace{\frac{\partial y_1^{in}}{\partial b_1^{(2)}}}_{= z_3 \cdot W_{31}^{(2)}} = \delta_1^{(2)}.$$



$$z = \varphi^{(1)}(W^{(1)}x + b^{(1)}),$$

$$y = \varphi^{(2)}(W^{(2)}z + b^{(2)}).$$

and $\frac{\partial E}{\partial b_2} = \delta_2^{(2)}$.

We see a pattern here ...

$$\frac{\partial E}{\partial W_{jk}^{(2)}} = \delta_k^{(2)} \cdot z_j, \quad \frac{\partial E}{\partial b_k} = \delta_k^{(2)}$$

/ Calculating $\delta_k^{(2)}$. . .

$$E(\theta) := \frac{1}{2} [(t_1 - y_1)^2 + (t_2 - y_2)^2]$$

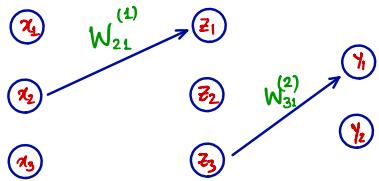
$$\delta_k^{(2)} = \frac{\partial E}{\partial y_k} \cdot \underbrace{\frac{\partial y_k}{\partial y_k^{in}}}_{=} = (y_k - t_k) \cdot [\varphi^{(2)}]'(y_k^{in}).$$

$$y_k = \varphi^{(2)}(y_k^{in}).$$

/ Similarly, let's calculate $\frac{\partial E}{\partial W_{12}^{(2)}}$.

$$\begin{aligned} \frac{\partial E}{\partial W_{12}^{(2)}} &= \frac{\partial E}{\partial y_2} \cdot \frac{\partial y_2}{\partial y_2^{in}} \cdot \frac{\partial y_2^{in}}{\partial W_{12}^{(2)}} \\ &= \delta_2^{(2)} \cdot z_1 \end{aligned}$$

" March on to calculate $\frac{\partial E}{\partial w_{jk}^{(1)}}$, $\frac{\partial E}{\partial b_k^{(1)}}$.



$$z = \varphi^{(1)}(w^{(1)}x + b^{(1)}),$$

$$y = \varphi^{(2)}(w^{(2)}z + b^{(2)}).$$

Let's tackle our example first. What is

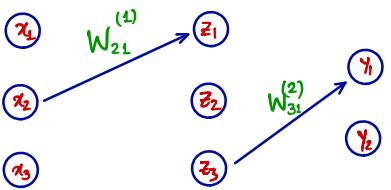
$$\frac{\partial E}{\partial w_{21}^{(1)}} ?$$

$$\frac{\partial E}{\partial w_{21}^{(1)}} = \underbrace{\frac{\partial E}{\partial z_1}}_{\text{Looks familiar}} \cdot \underbrace{\frac{\partial z_1}{\partial z_1^{\text{in}}}}_{\frac{\partial z_1^{\text{in}}}{\partial w_{21}^{(1)}}} \cdot \underbrace{\frac{\partial z_1^{\text{in}}}{\partial w_{21}^{(1)}}}_{\rightarrow x_2}.$$

Call this quantity $\delta_1^{(1)}$.

$$\therefore \frac{\partial E}{\partial w_{21}^{(1)}} = \delta_1^{(1)} \cdot x_2$$

... more generally $\frac{\partial E}{\partial w_{jk}^{(1)}} = \delta_k^{(1)} \cdot x_j$



$$z = \varphi^{(1)}(W^{(1)}x + b^{(1)}),$$

$$y = \varphi^{(2)}(W^{(2)}z + b^{(2)}).$$

$$\begin{aligned} \frac{\partial E}{\partial b_2^{(1)}} &= \frac{\partial E}{\partial y_2} \cdot \frac{\partial y_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2^{(1)}} \\ &= \delta_2^{(1)} \end{aligned}$$

... more generally

$$\frac{\partial E}{\partial b_k^{(l)}} = \delta_k^{(l)}.$$

" Calculating $\delta_k^{(l)}$.

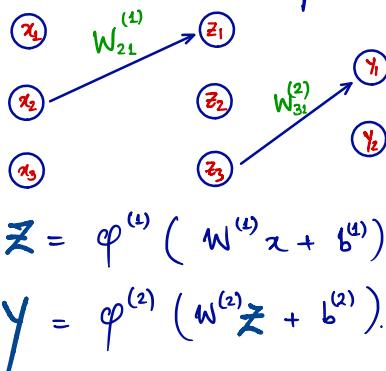
Notice that

$$\delta_k^{(l)} = \frac{\partial E}{\partial z_k} \cdot \frac{\partial z_k}{\partial z_k^{in}}$$

$$= \underbrace{\frac{\partial E}{\partial z_k}}_{\text{This quantity is not as easy to compute}} \cdot [\varphi^{(l)}]'(\underline{z_k^{in}}).$$

This quantity is not as easy to compute as $\frac{\partial E}{\partial y_k}$ in the last layer. Why? E depends explicitly on y_k 's and hence, it is easy to compute $\frac{\partial E}{\partial y_k}$, but E only depends implicitly on z 's.

Let's do a specific example: calculate $\frac{\partial E}{\partial y_1}$.



Key idea:
E depends on (z_1, z_2)
that in turn depends on
 y_1 .

$$\begin{aligned}\frac{\partial E}{\partial z_1} &= \frac{\partial E(y_1, y_2)}{\partial z_1} \\ &= \frac{\partial E}{\partial y_1} \cdot \frac{\partial y_1}{\partial z_1} + \frac{\partial E}{\partial y_2} \cdot \frac{\partial y_2}{\partial z_1} \\ &= \underbrace{\frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial y_1^{\text{in}}} \cdot \frac{\partial y_1^{\text{in}}}{\partial z_1}}_{= \delta_1^{(2)}} + \underbrace{\frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial y_2^{\text{in}}} \cdot \frac{\partial y_2^{\text{in}}}{\partial z_1}}_{\delta_2^{(2)} W_{12}^{(2)}}.\end{aligned}$$

$$\begin{aligned}y_1^{\text{in}} &= z_1 \cdot W_{11}^{(2)} + z_2 \cdot W_{21}^{(2)} + z_3 \cdot W_{31}^{(2)} \\ \Rightarrow \frac{\partial E}{\partial z_k} &= \delta_1^{(2)} W_{k1}^{(2)} + \delta_2^{(2)} W_{k2}^{(2)}.\end{aligned}$$

Notice that you
have already
calculated $\delta_1^{(2)}$ / $\delta_2^{(2)}$.

$$\therefore \delta_k^{(1)} = \frac{\partial E}{\partial z_k} \cdot [\varphi^{(1)}]'(z_k^{\text{in}})$$

$$= \left(\delta_1^{(2)} w_{k1}^{(2)} + \delta_2^{(2)} w_{k2} \right) [\varphi^{(1)}]'(z_k^{\text{in}}).$$

// Multi-layer neural network :

... left as an exercise.