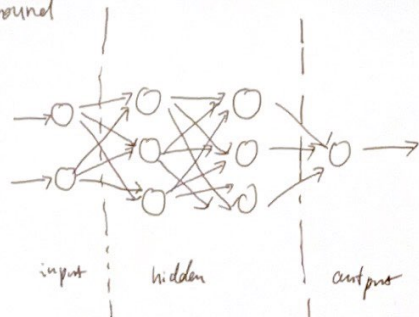


Background



$f_i()$: the i th fully connected layer

o : output

x : input ~~$x \in \mathbb{R}^{(b,1)}$~~

$$o = f_3(f_2(f_1(x)))$$

$$o = a(W_1x_1 + W_2x_2 + \dots + W_nx_n + b)$$

b : bias

a : activation function, a non linear function

w_i : the i th weight of the linear regression function

$$o^{(i)} = a(W^{(i)}x^{(i)} + b^{(i)})$$

$$o = a^{(3)}(W^{(3)}a^{(2)}(W^{(2)}a^{(1)}(W^{(1)}x^{(1)} + b^{(1)}) + b^{(2)}) + b^{(3)})$$

Stochastic gradient descent (SGD)

L : loss function

z : the output of the linear regression

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}$$

or

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial \theta}$$

to calculate the gradient of the weights

$$\text{chain rule: } \frac{d}{dx} f(g(h(x))) = \frac{d}{dx} f(g(h(x))) \cdot \frac{d}{dx} g(h(x)) \cdot \frac{d}{dx} h(x)$$

$$f'(g(h(x))) \cdot g'(h(x)) \cdot h'(x)$$

to calculate the gradient of the activation function

additive rule of derivation

$$\frac{d}{dx} (f(x) + g(x) + h(x)) = f'(x) + g'(x) + h'(x)$$

$$x \in \mathbb{R}^{(b,1)}$$

$$w \in \mathbb{R}^{(1,o)}$$

$$o_{\text{output}} \in \mathbb{R}^{(b,o)}$$

linear class

--init--

~~init~~ forward
backward
update_parameters

$$y = xW + b$$

update_parameters : • update the weights and biases of the network.

$$W_{ij} \leftarrow W_{ij} - \eta \frac{\partial L}{\partial W_{ij}}$$

$$\theta_j \leftarrow \theta_j - \eta \frac{\partial L}{\partial \theta_j}$$

* η : the learning rate

forward : • computes the regression against the inputs.

$$f = W_{1j}X_1 + W_{2j}X_2 + \dots + W_{nj}X_n + \theta_j$$

n : input size

W_{ij} : the weight for input entry i in regression unit j

θ_j : the bias term for regression unit j .

backward : given the gradient of the loss function with respect to the outputs

- computes the gradients of weights and biases.
- saves them as internal variable

$$\frac{\partial L}{\partial W_{ij}} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial W_{ij}} = O_i \cdot X_j \text{ or } \sum_k O_{ik} \cdot X_{jk}$$

$$\frac{\partial L}{\partial \theta_j} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial \theta_j} = O_j \text{ , or } \sum_k O_{jk}$$

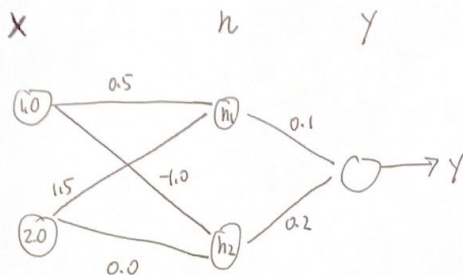
$$\frac{\partial L}{\partial X_j} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial X_j} = O_j \cdot W_j \text{ or } \sum_k O_{jk} \cdot W_{jk}$$

* i : index of the unit

j : index of the weight

k : index of the complete input

An Easy Example about how linear class works, especially the forward function



$$X = [1.0, 2.0]$$

$$W = \begin{bmatrix} 0.5 & -1.0 \\ 1.5 & 0.0 \end{bmatrix}$$

$$b = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} \begin{bmatrix} 0.1 & 0.2 \end{bmatrix}$$

$$y = X @ W + b$$

$$y = [1.0 \quad 2.0] \cdot \begin{bmatrix} 0.5 & -1.0 \\ 1.5 & 0.0 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} = [3.6 \quad -0.8]$$

$$X \in \mathbb{R}^{1 \times 2}$$

$$W \in \mathbb{R}^{2 \times 2}$$

$$b \in \mathbb{R}^{2 \times 1}$$

$$Y \in \mathbb{R}^{1 \times 2}$$

Generalize it

$$X \in \mathbb{R}^{1 \times n}$$

$$W \in \mathbb{R}^{n \times 2}$$

$$b \in \mathbb{R}^{n \times 1}$$

$$Y \in \mathbb{R}^{1 \times 2}$$

backwards

in forward propagation

$$z = x \cdot W + b$$

The Gradient of weights

$$\frac{\partial L}{\partial W} = \left(\frac{\partial z}{\partial W} \right)^T \cdot \frac{\partial L}{\partial z}$$

$$z = xW \Rightarrow \frac{\partial z}{\partial W} = x$$

$$\Rightarrow \frac{\partial L}{\partial W} = x^T \cdot \frac{\partial L}{\partial z}$$

e.g. $x = \begin{bmatrix} 1.0 & 2.0 \end{bmatrix}$

$$x^T = \begin{bmatrix} 1.0 \\ 2.0 \end{bmatrix}$$

$$\text{grad_output} = \begin{bmatrix} 1.0 & 2.0 \end{bmatrix}$$

$$\text{grad_W} = x^T @ \text{grad_output}$$

$$= \begin{bmatrix} 1.0 \\ 2.0 \end{bmatrix} \times \begin{bmatrix} 1.0 & 2.0 \end{bmatrix} = \begin{bmatrix} 1.0 & 2.0 \\ 2.0 & 4.0 \end{bmatrix}$$

L : loss function

W : weights Matrix

x : current input, can be batch

$\frac{\partial L}{\partial z}$: the gradient from last layer,
here called as grad_output

x^T : the transpose of x

The gradient of bias

$$\text{grad_biases} = \frac{\partial L}{\partial \theta_i} = \sum_k \frac{\partial L}{\partial \theta_i} O_{ik}$$

O_{ik} : grad_output[k,i]

in human words: the gradient of the
loss with respect to the output at dimension
 i for the k -th input sample.

$\frac{\partial L}{\partial \theta_i}$: The gradient of each bias θ_i is the sum
of the loss gradients at i -th output dimension
across all training samples.

e.g. grad_output = np.array([
 [0.1 -0.2]
 [0.3 -0.1]
 [0.2 0.0]
])

$$\text{np.sum(grad_output, axis=0, keepdims=True)}$$

$$= [0.1 + 0.3 + 0.2, -0.2 + (-0.1) + 0.0] = [0.6, 0.3]$$

axis aspect of matrix:

grad_output: $\begin{bmatrix} 0.1 & -0.2 \\ 0.3 & -0.1 \\ 0.2 & 0.0 \end{bmatrix}$ → sum up column wise

⑩ ~~Backward~~ Back propagation

The Gradient of output.

because: $z = xW + b$

$$\Rightarrow \frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \cdot W^T$$

$\frac{\partial L}{\partial z}$ means: How the loss L change if the output z changes?

it's the gradient, which flows into this layer from output.

it represents how the loss changes with the output of this layer, z

Therefore, $\frac{\partial L}{\partial z}$ is gradient_output

in python: $\text{grad_input} = \text{grad_output} @ \text{weight}^T$

an concret example:

$$X = \begin{bmatrix} 1 & 2 \end{bmatrix}$$

$$W = \begin{bmatrix} 0.5 & -1 \\ 1.5 & 0 \end{bmatrix}$$

$$b = \begin{bmatrix} 0.1 & 0.2 \end{bmatrix}$$

$$z = X @ W + b = \begin{bmatrix} 3.6 & -0.8 \end{bmatrix}$$

$$\text{grad_output} = dL/dz = \begin{bmatrix} 1 & 2 \end{bmatrix}$$

$$\text{grad_input} = \text{grad_output} @ W^T = \begin{bmatrix} -1.6 & 1.5 \end{bmatrix}$$