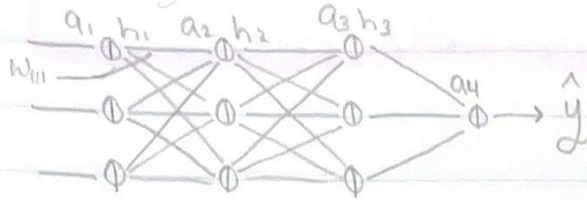


Gradient Descent

04/11/2023

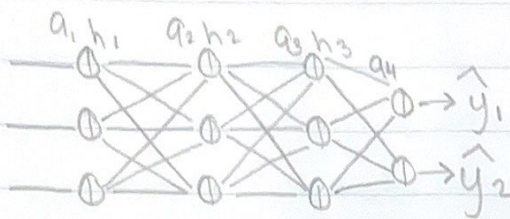
- defining the model



Regression

$$x \in \mathbb{R}^3$$

$$y \in \mathbb{R}$$



Classification

$$x \in \mathbb{R}^3$$

$$y \in \mathbb{R}^2$$

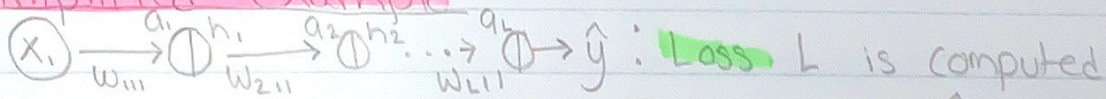
| | Real Values | Properties |
|-----------------|---------------|--|
| Output function | Linear layers | Softmax |
| Loss function | MSE | Cross Entropy = $-\sum_{i=1}^c p_i \log q_i$ |

w_{111} is updated by

$$w_{111} \leftarrow w_{111} - \alpha \frac{dL}{dw_{111}}$$

how do we compute this?

Simplified example



$$\frac{dL}{dw_{111}} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{da_1} \times \dots \times \frac{da_3}{dh_2} \times \frac{dh_2}{da_2} \times \frac{da_2}{da_1} \times \frac{da_1}{dw_{111}}$$

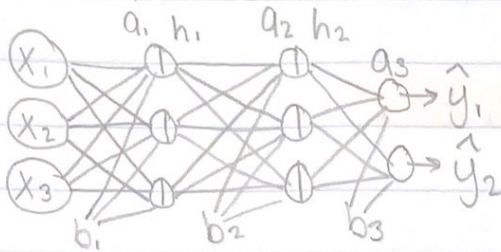
*only include w's that you want to optimize

$$\frac{dL}{dW_{ii}} = dL \frac{dh_i}{dW_{ii}}$$

(break chain between $\frac{da_2}{dh_i}$ and $\frac{dh_i}{da_1}$)

Gradient Descent (Part II)

04/13/2023



★ assuming a binary classification problem ★

Loss = L $\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$

★ $W_{iii} \leftarrow W_{iii} - \alpha \frac{dL}{dW_{iii}}$

$$\frac{dL}{dW_{iii}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{da_3} \cdot \frac{da_3}{dh_2} \cdot \frac{dh_2}{da_2} \cdot \frac{da_2}{dh_1} \cdot \frac{dh_1}{da_1} \cdot \frac{da_1}{dW_{iii}}$$

gradient across output layer

gradient from output layer to last hidden layer

gradient from 2nd hidden layer to 1st

gradient to the weight of interest

Quantities of interest

- 1) gradient with respect to output layer
- 2) gradient with respect to hidden layers
- 3) gradient with respect to weights and biases

For classification:

cost function: cross entropy

output layer: soft max

example (cross entropy):

$\hat{y} = \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}$ true prob = $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$CE = -\sum p_i \log q_i = -\log(0.8)$

★ Negative log likelihood ★

parameters (for example):

$$\theta = [w_1, w_2, w_3, b_1, b_2, b_3]$$

$$\textcircled{1} L = -\log(\hat{y}_l) \quad \{l = \text{true class label}$$

$$\frac{dL}{d\hat{y}_i} = \begin{cases} -\frac{1}{\hat{y}_i} & \text{if } i=l \\ 0 & \text{otherwise} \end{cases} = -\frac{\mathbb{1}_{i=l}}{\hat{y}_i} \quad \left\{ \mathbb{1}_{i=l} = \begin{cases} 1 & i=l \\ 0 & \text{otherwise} \end{cases} \right\}$$

indicator function

$$\nabla_{\hat{y}} L = \begin{bmatrix} \frac{dL}{d\hat{y}_1} \\ \frac{dL}{d\hat{y}_2} \\ \vdots \end{bmatrix} = -\frac{1}{\hat{y}_l} \begin{bmatrix} \mathbb{1}_{l=1} \\ \mathbb{1}_{l=2} \\ \vdots \end{bmatrix} = -\frac{1}{\hat{y}_l} e(l)$$

is 1 the true class?

where $e(l)$ = a vector of all zeros with length = $\text{len}(\hat{y})$, and a one at l^{th} position

$$\nabla_{\hat{y}} L = -\frac{1}{\hat{y}_l} e(l)$$

$$\frac{d\hat{y}}{da_{\text{last}}} \quad \left(\text{in our example, this is } \frac{d\hat{y}}{da_3} \right)$$

simplest: $\frac{d\hat{y}_l}{da_{\text{last } i}} = \frac{d}{da_{\text{last } i}} \left(\frac{\exp(a_{\text{last } l})}{\sum_j \exp(a_{\text{last } j})} \right)$

$$\frac{d\hat{y}_l}{da_{\text{last } i}} = \hat{y}_l (\mathbb{1}_{l=i} - \hat{y}_i)$$

$$\frac{d\hat{y}}{da_{\text{last } i}} = \hat{y} (\mathbb{1}_{l=i} - \hat{y}_i)$$

$$\frac{d\hat{y}}{d\hat{a}_{last\ i}} = \begin{bmatrix} \frac{dy_1}{d\hat{a}_{last\ i}} \\ \frac{dy_2}{d\hat{a}_{last\ i}} \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \hat{y}_1 (\mathbb{1}_{1=i} - \hat{y}_i) \\ \hat{y}_2 (\mathbb{1}_{2=i} - \hat{y}_i) \\ \vdots \\ \vdots \end{bmatrix} = \hat{y} (\mathbb{1}_{l=i} - \hat{y}_i)$$

$$\frac{dL}{d\hat{a}_{last}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\hat{a}_{last}}$$

$$\begin{aligned} \frac{dL}{d\hat{a}_{last\ i}} &= \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{d\hat{a}_{last\ i}} \\ &= \frac{-1}{\hat{y}_e} e(l) \cdot \hat{y} (\mathbb{1}_{l=i} - \hat{y}_i) \\ &= -(\mathbb{1}_{l=i} - \hat{y}_i) \end{aligned}$$

$$\frac{dL}{d\hat{a}_{last}} = \begin{bmatrix} \frac{dL}{d\hat{a}_{last\ 1}} \\ \frac{dL}{d\hat{a}_{last\ 2}} \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} -(\mathbb{1}_{l=1} - \hat{y}_1) \\ -(\mathbb{1}_{l=2} - \hat{y}_2) \\ \vdots \\ \vdots \end{bmatrix} = -(e(l) - \hat{y})$$

$$\boxed{\frac{dL}{d\hat{a}_{last}} = -(e(l) - \hat{y})}$$