

Lecture Notes:

Building a Neural Network From Scratch

Introduction:

- I'll assume an understanding of fundamental math concepts, like multivariable calculus & linear algebra, as well as a general understanding of the structure of neural networks, and programming.
- What I want to focus on is the precise mathematical formulation of a simple feedforward neural network, and then translate that into an actual model in Python.
- If you're shaky on some of the fundamental concepts, that's OK - you'll build an intuition moving forward.

Talk about problems

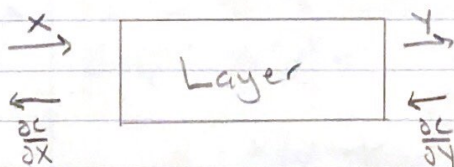
Implementation Architecture:

Steps in Machine Learning:

- ① Data \rightarrow Model
- ② Compute the distance between desired and actual output
- ③ Adjust parameters of model
 \rightarrow Repeat

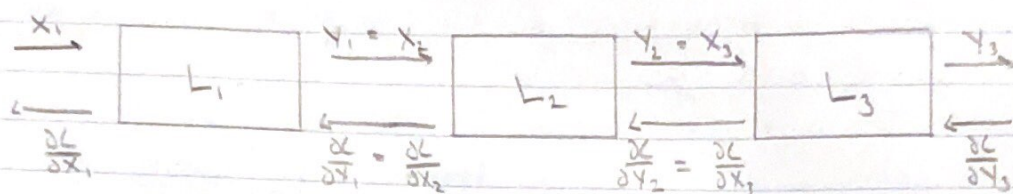
Before starting anything, look at data
 \downarrow

- Goal is modular code \Rightarrow implement every layer separately



- With this design, we can think of each layer as an individual object

- Sequential Model: Output of one layer is input of the next



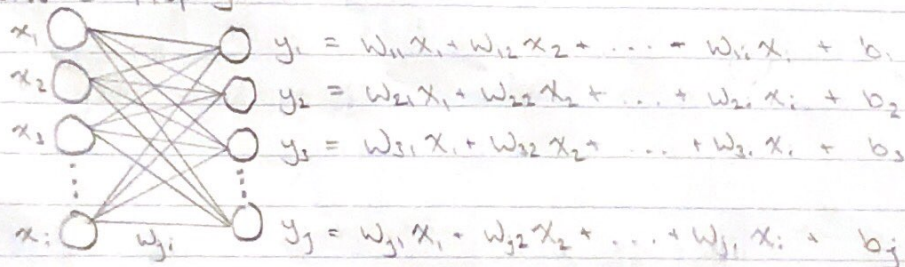
- This describes our Forward propagation and backward propagation algorithms

→ Implement Base Layer

Dense Layer: aka "Fully Connected"

- Each input node is connected to each output node

Forward Propagation:



⇒ Matrix Equation: $Y = WX + b$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix}_{j \times 1} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1i} \\ w_{21} & w_{22} & \dots & w_{2i} \\ \vdots & \vdots & \ddots & \vdots \\ w_{j1} & w_{j2} & \dots & w_{ji} \end{bmatrix}_{j \times i} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_i \end{bmatrix}_{i \times 1} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_j \end{bmatrix}_{j \times 1}$$

Back propagation:

Given a cost function C , we want to know how to change the parameters to reduce the error

⇒ Gradient Descent

Compute $\frac{\partial C}{\partial y}$, we want $\frac{\partial C}{\partial w}$, $\frac{\partial C}{\partial b}$, and $\frac{\partial C}{\partial x}$

$$\frac{\partial C}{\partial y} = \begin{bmatrix} \frac{\partial C}{\partial y_1} \\ \frac{\partial C}{\partial y_2} \\ \vdots \\ \frac{\partial C}{\partial y_j} \end{bmatrix}_{j \times 1} \quad \xrightarrow{\text{We want:}} \quad \frac{\partial C}{\partial w} = \begin{bmatrix} \frac{\partial C}{\partial w_{11}} & \frac{\partial C}{\partial w_{12}} & \dots & \frac{\partial C}{\partial w_{1n}} \\ \frac{\partial C}{\partial w_{21}} & \frac{\partial C}{\partial w_{22}} & \dots & \frac{\partial C}{\partial w_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial C}{\partial w_{j1}} & \frac{\partial C}{\partial w_{j2}} & \dots & \frac{\partial C}{\partial w_{jn}} \end{bmatrix}_{j \times i}$$

Using chain rule: $\frac{\partial C}{\partial w} = \frac{\partial C}{\partial y} \cdot \frac{\partial y}{\partial w}$

$$\text{e.g. } \frac{\partial C}{\partial w_{12}} = \frac{\partial C}{\partial y_1} \underbrace{\frac{\partial y_1}{\partial w_{12}}}_{=x_1} + \underbrace{\frac{\partial C}{\partial y_2} \frac{\partial y_2}{\partial w_{12}} + \dots + \frac{\partial C}{\partial y_j} \frac{\partial y_j}{\partial w_{12}}}_{\text{these} = 0!}$$

$$\Rightarrow \frac{\partial C}{\partial w_{ji}} = \frac{\partial C}{\partial y_j} x_i$$

$$\text{So, } \frac{\partial C}{\partial w} = \begin{bmatrix} \frac{\partial C}{\partial y_1} x_1 & \frac{\partial C}{\partial y_1} x_2 & \dots & \frac{\partial C}{\partial y_1} x_i \\ \frac{\partial C}{\partial y_2} x_1 & \frac{\partial C}{\partial y_2} x_2 & \dots & \frac{\partial C}{\partial y_2} x_i \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial C}{\partial y_j} x_1 & \frac{\partial C}{\partial y_j} x_2 & \dots & \frac{\partial C}{\partial y_j} x_i \end{bmatrix}_{j \times i} = \begin{bmatrix} \frac{\partial C}{\partial y} \\ \vdots \\ \frac{\partial C}{\partial y_j} \end{bmatrix}_{j \times 1} X^T_{1 \times i}$$

$$\frac{\partial L}{\partial \mathbf{b}} = \begin{bmatrix} \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial b_2} \\ \vdots \\ \frac{\partial L}{\partial b_j} \end{bmatrix}_{j \times 1}$$

$$\text{ex) } \frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial y_1} \underbrace{\frac{\partial y_1}{\partial b_1}}_{=1} + \underbrace{\frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial b_1} + \dots + \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial b_1}}_{=0}$$

$$\text{Thus, } \frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}}$$

Finally, we need $\frac{\partial L}{\partial \mathbf{x}}$ to pass back to previous input

$$\frac{\partial L}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial L}{\partial x_1} \\ \frac{\partial L}{\partial x_2} \\ \vdots \\ \frac{\partial L}{\partial x_i} \end{bmatrix}_{i \times 1}$$

$$\begin{aligned} \text{ex) } \frac{\partial L}{\partial x_1} &= \frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial x_1} + \frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial x_1} + \dots + \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial x_1} \\ &= \frac{\partial L}{\partial y_1} w_{11} + \frac{\partial L}{\partial y_2} w_{21} + \dots + \frac{\partial L}{\partial y_j} w_{j1} \end{aligned}$$

$$\text{So, } \frac{\partial L}{\partial \mathbf{x}} = \mathbf{W}^T \frac{\partial L}{\partial \mathbf{y}}$$

Activation Layer:

Computing our activation Function element-wise

$$x_1 \rightarrow \bigcirc \rightarrow y_1 = F(\cdot)$$

$$x_2 \rightarrow \bigcirc \rightarrow y_2 = F(\cdot)$$

$$x_3 \rightarrow \bigcirc \rightarrow y_3 = F(\cdot) \Rightarrow Y = F(X)$$

\vdots

$$x_i \rightarrow \bigcirc \rightarrow y_i = F(\cdot)$$

$$\frac{\partial L}{\partial Y} = \begin{bmatrix} \frac{\partial L}{\partial y_1} \\ \frac{\partial L}{\partial y_2} \\ \vdots \\ \frac{\partial L}{\partial y_i} \end{bmatrix}_{i \times 1}$$

We want \rightarrow

$$\frac{\partial L}{\partial X} = \begin{bmatrix} \frac{\partial L}{\partial x_1} \\ \frac{\partial L}{\partial x_2} \\ \vdots \\ \frac{\partial L}{\partial x_i} \end{bmatrix}_{i \times 1}$$

$$\text{ex) } \frac{\partial L}{\partial x_1} = \underbrace{\frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial x_1}}_{=0} + \underbrace{\frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial x_1}}_{=0} + \dots + \underbrace{\frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial x_1}}_{=0}$$

$$\Rightarrow \frac{\partial L}{\partial x_1} = \frac{\partial L}{\partial y_1} \cdot F'(x_1); \quad \text{So, } \frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \odot F'(X)$$

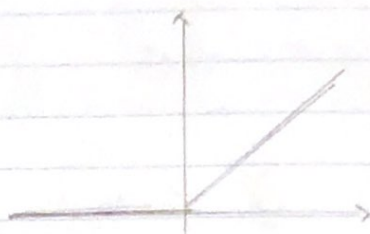
\uparrow
Hadamard Product,

i.e., element-wise multiplication

* Note: We need a nonlinear activation Function, otherwise our network just computes a linear combination of the input.

ReLU: $g(x) = \max(0, x)$

$$g'(x) = \begin{cases} 0, & x \leq 0 \\ 1, & \text{else} \end{cases}$$



Softmax: $y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$

$$\frac{\partial \mathcal{L}}{\partial x_k} = \frac{\partial \mathcal{L}}{\partial y_1} \frac{\partial y_1}{\partial x_k} + \frac{\partial \mathcal{L}}{\partial y_2} \frac{\partial y_2}{\partial x_k} + \dots + \frac{\partial \mathcal{L}}{\partial y_n} \frac{\partial y_n}{\partial x_k}$$

If $k=i$: $\frac{\partial y_i}{\partial x_k} = \frac{e^{x_i} (\sum_{j=1}^n e^{x_j}) - e^{x_i} e^{x_i}}{(\sum_{j=1}^n e^{x_j})^2}$ ← Quotient Rule

$$= \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} - \left(\frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \right)^2$$

$$= y_i (1 - y_i)$$

If $k \neq j$: $\frac{\partial y_i}{\partial x_k} = \frac{-e^{x_k} e^{x_i}}{(\sum_{j=1}^n e^{x_j})^2} = -y_k y_i$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial x} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_1} \frac{\partial y_1}{\partial x_1} + \dots + \frac{\partial \mathcal{L}}{\partial y_n} \frac{\partial y_n}{\partial x_1} \\ \frac{\partial \mathcal{L}}{\partial y_1} \frac{\partial y_1}{\partial x_2} + \dots + \frac{\partial \mathcal{L}}{\partial y_n} \frac{\partial y_n}{\partial x_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial y_1} \frac{\partial y_1}{\partial x_n} + \dots + \frac{\partial \mathcal{L}}{\partial y_n} \frac{\partial y_n}{\partial x_n} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_n} \end{bmatrix} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_1} \\ \frac{\partial \mathcal{L}}{\partial y_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial y_n} \end{bmatrix}$$

$$= \begin{bmatrix} y_1(1-y_1) & -y_1y_2 & \dots & -y_1y_n \\ -y_2y_1 & y_2(1-y_2) & \dots & -y_2y_n \\ \vdots & \vdots & \ddots & \vdots \\ -y_ny_1 & -y_ny_2 & \dots & y_n(1-y_n) \end{bmatrix} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_1} \\ \frac{\partial \mathcal{L}}{\partial y_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial y_n} \end{bmatrix}$$

$$= \left(\begin{bmatrix} y_1 & y_2 & \dots & y_n \\ y_2 & y_2 & \dots & y_2 \\ \vdots & \vdots & \ddots & \vdots \\ y_n & y_n & \dots & y_n \end{bmatrix} \odot \begin{bmatrix} 1-y_1 & -y_2 & \dots & -y_n \\ -y_1 & 1-y_2 & \dots & -y_n \\ \vdots & \vdots & \ddots & \vdots \\ -y_1 & -y_2 & \dots & 1-y_n \end{bmatrix} \right) \cdot \frac{\partial \mathcal{L}}{\partial y}$$

$$= \left(\begin{bmatrix} y_1 & y_2 & \dots & y_n \\ y_2 & y_2 & \dots & y_2 \\ \vdots & \vdots & \ddots & \vdots \\ y_n & y_n & \dots & y_n \end{bmatrix} \odot \left(I - \begin{bmatrix} y_1 & y_2 & \dots & y_n \\ y_1 & y_2 & \dots & y_n \\ \vdots & \vdots & \ddots & \vdots \\ y_1 & y_2 & \dots & y_n \end{bmatrix} \right) \right) \cdot \frac{\partial \mathcal{L}}{\partial y}$$

$\hookrightarrow A$

$$= (A \odot (I - A^T)) \cdot \frac{\partial \mathcal{L}}{\partial y}$$

Cost Function: Mean-Squared Error

Given desired output Y^* , actual output Y ,

$$Y^* = \begin{bmatrix} y_1^* \\ y_2^* \\ \vdots \\ y_n^* \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix},$$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i^* - y_i)^2$$

$$\text{ex) } \frac{\partial C}{\partial y_1} = \frac{1}{n} \frac{\partial}{\partial y_1} \left[(y_1^* - y_1)^2 + \underbrace{(y_2^* - y_2)^2 + \dots + (y_n^* - y_n)^2}_{=0} \right]$$

$$= \frac{2}{n} (y_1 - y_1^*)$$

$$\Rightarrow \frac{\partial C}{\partial Y} = \frac{2}{n} (Y - Y^*)$$