

Derivation of Backpropagation

Introduction (Basic derivation)

$$\text{Let } \hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

and $\text{Loss} = (y - \hat{y})^2 \leftarrow \text{squared error loss}$

we can rewrite the above in "matrix" format.

$$\underset{1 \times 1}{[y]} = \underset{1 \times 3}{[x_1 \ x_2 \ x_3]} \underset{3 \times 1}{\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}} + \underset{1 \times 1}{[b]}$$

What we want is to learn the best W and b , and

for backprop we need derivative of learnable parameter

w.r.t. Loss.

Notice gradients have same shape as the original tensors!

$$\text{So, } \frac{dL}{dW} = \underset{3 \times 1}{\begin{bmatrix} dL/dw_1 \\ dL/dw_2 \\ dL/dw_3 \end{bmatrix}} \quad \frac{dL}{db} = \underset{1 \times 1}{[dL/db]}$$

Lets do dL/dw_1 first!

$$\frac{dL}{dw_1} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_1}$$

chain rule

$$\text{if } L = (y - \hat{y})^2, \frac{dL}{d\hat{y}} = -2(y - \hat{y})$$

$$\text{if } \hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b, \frac{d\hat{y}}{dw_1} = x_1$$

$$\therefore \frac{dL}{dw_1} = -2(y - \hat{y})x_1$$

$$\text{More generally, } \frac{dL}{d\mathbf{w}} = \begin{bmatrix} dL/dw_1 \\ dL/dw_2 \\ dL/dw_3 \end{bmatrix} = \begin{bmatrix} -2(y - \hat{y})x_1 \\ -2(y - \hat{y})x_2 \\ -2(y - \hat{y})x_3 \end{bmatrix}$$

3×1

$$= \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \begin{bmatrix} -2(y - \hat{y}) \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}^T \frac{dL}{d\hat{y}}$$

3×1 1×1

Remember! \vec{X} was a row vector, now its a column vector!

Now we need our gradient for the intercept term!

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{db}$$

$$\therefore \frac{dL}{db} = \frac{dL}{d\hat{y}} \cdot 1$$

$$\frac{dL}{d\hat{y}} = -2(y - \hat{y})$$

(same as before)

$$\frac{d\hat{y}}{db} = 1$$

$$\text{Final Results: } \frac{dL}{dw} = \vec{x}^T \frac{dL}{d\hat{y}}$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}}$$

Adding a Batch Dimension

A thing we omitted in the past derivation was the batch dimension! In neural networks we typically pass in N samples and then propagate back to average loss!

So let's rewrite the problem to have a batch size of 2.

Let:

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$\oplus \begin{bmatrix} b \end{bmatrix}_{1 \times 1}$$



Element-wise sum!
same $b_{1,1}$ is added to
each sample in the batch
independently (broadcasting)

And our loss is now Mean Squared Error

$$L = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad \text{where } N=2 \text{ in our case.}$$

just line before, we need to compute $\frac{dL}{d\hat{y}}$, but we

will compute this for all \hat{y}_i as we have multiple samples in our batch.

move constant into sum.

$$\frac{dL}{d\hat{y}} \left(\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right)$$

Derivative of a sum is the sum of the derivatives.

$$= \frac{dL}{d\hat{y}} \sum_{i=1}^N \frac{1}{N} (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^N \frac{dL}{d\hat{y}_i} \frac{1}{N} (y_i - \hat{y}_i)^2$$

$$\frac{dL}{d\hat{y}} = \sum_{i=1}^N -\frac{2}{N} (y_i - \hat{y}_i)$$

↑
We haven't dealt w/ this sum yet, but we will soon!

easy to compute derivative same as before!

So now, just as before,

$$\frac{dL}{dW} = \begin{bmatrix} dL/dw_1 \\ dL/dw_2 \\ dL/dw_3 \end{bmatrix}$$

Let's do dL/dw_1 first.

$$\frac{dL}{dw_1} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_1} = \sum_{i=1}^2 -\frac{2}{2} (y_i - \hat{y}_i) \frac{d\hat{y}_i}{dw_1}$$

how do we differentiate \hat{y}_i wrt w_1 ?

$$\hat{y}_1 = w_1 x_{11} + w_2 x_{12} + w_3 x_{13} + b$$

$$\hat{y}_2 = w_1 x_{21} + w_2 x_{22} + w_3 x_{23} + b$$

$$\text{then, } \frac{d\hat{y}_1}{dw_1} = x_{11} \quad \text{and} \quad \frac{d\hat{y}_2}{dw_1} = x_{21}$$

so then

$$\frac{dL}{dw_1} = \sum_{i=1}^2 -\frac{2}{2} (y_i - \hat{y}_i) \frac{d\hat{y}_i}{dw_1}$$

$$= -(y_1 - \hat{y}_1) x_{11} - (y_2 - \hat{y}_2) x_{21}$$

Similarly:

$$\frac{dL}{dw_2} = -(y_1 - \hat{y}_1)X_{12} - (y_2 - \hat{y}_2)X_{22}$$

$$\frac{dL}{dw_3} = -(y_1 - \hat{y}_1)X_{13} - (y_2 - \hat{y}_2)X_{23}$$

All together!

Sum over batch is handled within the matrix multiplication

$$\frac{dL}{dW} = \begin{bmatrix} dL/dw_1 \\ dL/dw_2 \\ dL/dw_3 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} -(y_1 - \hat{y}_1)X_{11} - (y_2 - \hat{y}_2)X_{21} \\ -(y_1 - \hat{y}_1)X_{12} - (y_2 - \hat{y}_2)X_{22} \\ -(y_1 - \hat{y}_1)X_{13} - (y_2 - \hat{y}_2)X_{23} \end{bmatrix}_{3 \times 1}$$

$$= \begin{bmatrix} X_{11} & X_{21} \\ X_{12} & X_{22} \\ X_{13} & X_{23} \end{bmatrix} \begin{bmatrix} -(y_1 - \hat{y}_1) \\ -(y_2 - \hat{y}_2) \end{bmatrix}$$

$$= X^T \begin{bmatrix} dL/d\hat{y}_1 \\ dL/d\hat{y}_2 \end{bmatrix} = X^T \frac{dL}{d\hat{y}}$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \cdot \underbrace{\frac{d\hat{y}}{db}}_{\text{always 1}} = \sum_{i=1}^2 -(y_i - \hat{y}_i) \frac{d\hat{y}_i}{db} \rightarrow 1$$

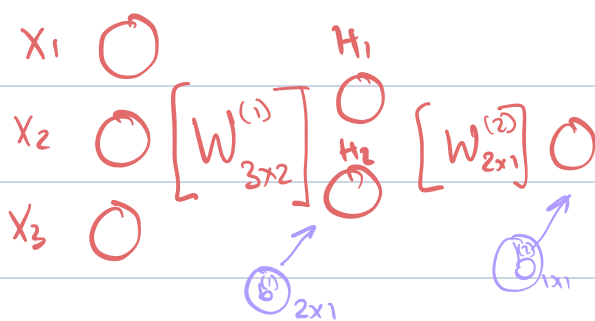
$$= \sum_{i=1}^2 \frac{dL}{d\hat{y}_i}$$

$$\therefore \left[\frac{dL}{dw} = X^T \frac{dL}{d\hat{y}} \quad \frac{dL}{db} = \sum_{i=1}^N \frac{dL}{d\hat{y}_i} \right]$$

Stacking Layers!

Until now we have done a very simple projection with a Weight W and bias b .

But deep learning stacks multiple layers, what does that look like?



$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix} \quad H = \begin{bmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{bmatrix}$$

$$W^{(1)} = \begin{bmatrix} W_{3 \times 2} \end{bmatrix}$$

$$W^{(2)} = \begin{bmatrix} W_{2 \times 1} \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} b_{2 \times 1} \end{bmatrix}$$

$$b^{(2)} = \begin{bmatrix} b_{1 \times 1} \end{bmatrix}$$

$$\hat{y} = [XW^{(1)} + b^{(1)}]W^{(2)} + b^{(2)}$$

$$Loss = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

So split this into 2 composed functions:

$$\hat{y} = HW^{(2)} + b^{(2)} \quad H = XW^{(1)} + b^{(1)}$$

So let's find the gradients of $W^{(2)}$ and $b^{(2)}$ based on what we did previously:

$$\frac{dL}{dW^{(2)}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dW^{(2)}} = H^T \frac{dL}{d\hat{y}}$$

$$\frac{dL}{db^{(2)}} = \sum_{i=1}^N \frac{dL}{d\hat{y}_i}$$

Now that we have gradients for $W^{(2)}$ and $b^{(2)}$

we must backpropagate to $W^{(1)}$ and $b^{(1)}$

$$\frac{dL}{dW^{(1)}} = \frac{dL}{dH} \cdot \frac{dH}{dW^{(1)}} = X^T \frac{dL}{dH} \quad \leftarrow \text{Using previous formula}$$

$$\frac{dL}{dH} = \underbrace{\frac{dL}{d\hat{y}}}_{\text{we know this!}} \cdot \underbrace{\frac{d\hat{y}}{dH}}_{\text{we don't know this!}}$$

we have not computed this intermediate, but we need it for backprop to a new layer!

Remember, the H in our case is
we have a batch size of 2

$$\begin{bmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{bmatrix} \quad b/c$$

$B \times 2$

$$\frac{d\hat{y}}{dH} = \begin{bmatrix} dL/dh_{11} & dL/dh_{12} \\ dL/dh_{21} & dL/dh_{22} \end{bmatrix}$$

$$\text{Let } \frac{dL}{dh_{11}} = \frac{dL}{d\hat{y}_1} \cdot \frac{d\hat{y}_1}{dh_{11}}$$

$$\text{if } \hat{y}_1 = w_1^{(2)} h_{11} + w_2^{(2)} h_{12} + b^{(2)}$$

$$\text{Then } \frac{d\hat{y}_1}{dh_{11}} = w_1^{(2)}$$

$$\text{Similarly, } \frac{d\hat{y}_1}{dh_{12}} = w_2^{(2)} \quad \frac{d\hat{y}_2}{dh_{21}} = w_1^{(1)} \quad \text{and} \quad \frac{d\hat{y}_2}{dh_{22}} = w_2^{(1)}$$

$$\text{Then } \frac{d\hat{y}}{dH} = \begin{bmatrix} d\hat{y}_1/dh_{11} & d\hat{y}_1/dh_{12} \\ d\hat{y}_2/dh_{21} & d\hat{y}_2/dh_{22} \end{bmatrix} = \begin{bmatrix} w_1^{(2)} & w_2^{(2)} \\ w_1^{(1)} & w_2^{(1)} \end{bmatrix}$$

$$\frac{dL}{dH} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dH} = \begin{bmatrix} \frac{dL}{d\hat{y}_1} \cdot \frac{d\hat{y}_1}{dh_{11}} & \frac{dL}{d\hat{y}_1} \cdot \frac{d\hat{y}_1}{dh_{12}} \\ \frac{dL}{d\hat{y}_2} \cdot \frac{d\hat{y}_2}{dh_{21}} & \frac{dL}{d\hat{y}_2} \cdot \frac{d\hat{y}_2}{dh_{22}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{dL}{d\hat{y}_1} \cdot w_1^{(2)} & \frac{dL}{d\hat{y}_1} \cdot w_2^{(2)} \\ \frac{dL}{d\hat{y}_2} \cdot w_1^{(1)} & \frac{dL}{d\hat{y}_2} \cdot w_2^{(1)} \end{bmatrix} = \begin{bmatrix} dL/d\hat{y}_1 \\ dL/d\hat{y}_2 \end{bmatrix} \begin{bmatrix} w_1^{(2)} & w_2^{(2)} \\ w_1^{(1)} & w_2^{(1)} \end{bmatrix}$$

$2 \times 1 \quad 1 \times 2$

$$= \frac{dL}{d\hat{y}} \cdot W^T$$

$$\therefore \frac{dL}{dw^{(1)}} = \underbrace{\frac{dL}{d\hat{y}}}_{\text{known}} \cdot \frac{d\hat{y}}{dH} \cdot \underbrace{\frac{dH}{dw^{(1)}}}_{\text{we know the formula}}$$

$$\therefore \frac{dL}{dw^{(1)}} = X^T \frac{dL}{d\hat{y}} W^T$$

Similarly, $\frac{dL}{db^{(1)}} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dH} \cdot \frac{dH}{db^{(1)}}$ → always 1

$$\frac{dL}{db^{(1)}} = \sum_{i=1}^N \frac{dL}{d\hat{y}_i} \cdot \frac{d\hat{y}_i}{dH}$$

So all the formulas we need are:

if $\hat{y} = XW + b$ and L is some loss.

$$\boxed{\frac{dL}{dw} = X^T \frac{dL}{d\hat{y}} \quad \frac{dL}{db} = \sum_{i=1}^N \frac{dL}{d\hat{y}_i}}$$

and for backprop on multiple stacked layers,

we need $\boxed{\frac{dL}{dX} = \frac{dL}{d\hat{y}} W^T}$

Adding an Activation Function

All we did so far is stack 2 linear layers, but what if we also add an activation function?

$$\hat{y} = \sigma(XW^{(1)} + b^{(1)})W^{(2)} + b^{(2)}$$

So our composition of functions now is:

$$H = XW^{(1)} + b^{(1)}$$

$$S = \sigma(H)$$

$$\hat{y} = SW^{(2)} + b^{(2)}$$

First, find gradients on $W^{(2)}$ and $b^{(2)}$

$$\frac{dL}{dW^{(2)}} = S^T \frac{dL}{d\hat{y}} \quad \frac{dL}{db^{(2)}} = \sum_{i=1}^N 1 \frac{dL}{d\hat{y}_i}$$

$$\frac{dL}{dS} = \frac{dL}{d\hat{y}} \cdot W^{(2)T}$$

now we need to backpropagate but through the activation function.

new derivative for sigmoid

$$\frac{dL}{dw^{(1)}} = \frac{dL}{ds} \cdot \frac{ds}{dh} \cdot \frac{dh}{dw^{(1)}}$$

$$= \frac{dL}{d\hat{y}} w^{(2)T} \left(\underbrace{\sigma(h)(1-\sigma(h))}_{(\sigma(h))' = \sigma(h)(1-\sigma(h))} \right) \frac{dh}{dw^{(1)}}$$

$$= X^T \frac{dL}{d\hat{y}} w^{(2)T} (\sigma(h)(1-\sigma(h)))$$

$$\frac{dL}{db^{(2)}} = \sum_{i=1}^N \frac{dL}{d\hat{y}} w^{(2)T} (\sigma(h)(1-\sigma(h)))$$

Summary!

At every linear layer compute

$\frac{dL}{dw}$, $\frac{dL}{db}$ and use them to update weights and intercept.

also compute $\frac{dL}{dx}$ as it is needed for chain rule for more previous layers.