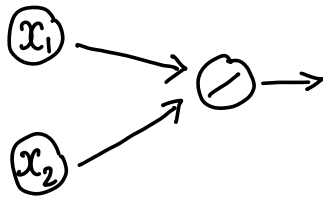
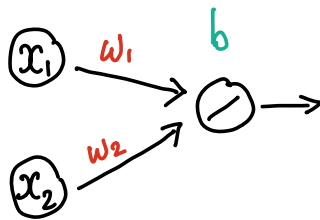


Consider a two-input problem with a continuous dependent variable and no hidden layers.



What are the parameters to be learned?



The model
and MSE
loss

$$\begin{aligned}\text{model}(x_1, x_2) &= b + w_1 x_1 + w_2 x_2 \\ \text{loss} &= (\text{model}(x_1, x_2) - y)^2 \\ \text{loss} &= (b + w_1 x_1 + w_2 x_2 - y)^2\end{aligned}$$

Let's calculate $\nabla \text{loss} = \left[\frac{\partial \text{loss}}{\partial b} \quad \frac{\partial \text{loss}}{\partial w_1} \quad \frac{\partial \text{loss}}{\partial w_2} \right]$ the "old-fashioned" way 😊

$$\begin{aligned} \text{loss} &= (b + w_1 x_1 + w_2 x_2 - y)^2 \\ &\begin{cases} \frac{\partial \text{loss}}{\partial b} = 2(b + w_1 x_1 + w_2 x_2 - y) \\ \frac{\partial \text{loss}}{\partial w_1} = 2(b + w_1 x_1 + w_2 x_2 - y) x_1 \\ \frac{\partial \text{loss}}{\partial w_2} = 2(b + w_1 x_1 + w_2 x_2 - y) x_2 \end{cases} \end{aligned}$$

Now, let's organize the calculations a bit differently

$$\text{Let } a_1 = w_1 x_1$$

$$a_2 = w_2 x_2$$

$$\hat{y} = b + a_1 + a_2$$

Plugging \nearrow into \downarrow

$$\text{Loss} = (b + w_1 x_1 + w_2 x_2 - y)^2$$

we get:

$$\text{Loss} = (\hat{y} - y)^2$$

$$\frac{\partial \text{Loss}}{\partial b} = ?$$

$$\frac{\partial \text{Loss}}{\partial w_1} = ?$$

$$\frac{\partial \text{Loss}}{\partial w_2} = ?$$

We can
apply the
Chain Rule!

$$\begin{aligned} \text{Loss} &= (\hat{y} - y)^2 \\ \hat{y} &= b + a_1 + a_2 \\ a_1 &= w_1 x_1 \\ a_2 &= w_2 x_2 \end{aligned} \quad \Rightarrow \quad \begin{aligned} \frac{\partial \text{Loss}}{\partial b} &= \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} \\ \frac{\partial \text{Loss}}{\partial w_1} &= \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} \\ \frac{\partial \text{Loss}}{\partial w_2} &= \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_2} \end{aligned}$$

Next, we need to calculate

$$\frac{\partial \text{Loss}}{\partial \hat{y}}, \quad \frac{\partial \hat{y}}{\partial b}, \quad \frac{\partial a_1}{\partial w_1} \text{ and } \frac{\partial a_2}{\partial w_2}$$

But that's easy!

$$\text{Loss} = (\hat{y} - y)^2$$



$$\frac{\partial \text{Loss}}{\partial \hat{y}} = 2(\hat{y} - y)$$

$$\hat{y} = b + a_1 + a_2$$



$$\frac{\partial \hat{y}}{\partial b} = \frac{\partial \hat{y}}{\partial a_1} = \frac{\partial \hat{y}}{\partial a_2} = 1$$

$$a_1 = w_1 x_1$$



$$\frac{\partial a_1}{\partial w_1} = x_1$$

$$a_2 = w_2 x_2$$



$$\frac{\partial a_2}{\partial w_2} = x_2$$

Putting everything together:

$$\frac{\partial \text{Loss}}{\partial b} = \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} = 2(\hat{y} - y) \cdot 1$$

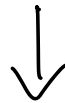
$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} = 2(\hat{y} - y) \cdot 1 \cdot x_1$$

$$\frac{\partial \text{Loss}}{\partial w_2} = \frac{\partial \text{Loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_2} = 2(\hat{y} - y) \cdot 1 \cdot x_2$$

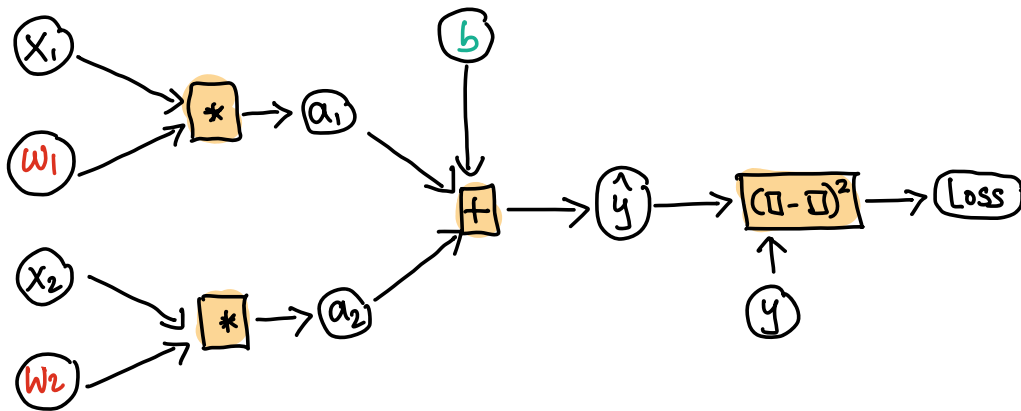
You can check this matches the "old fashioned" calculation!!

OK, we are finally ready for backpropagation!! 😊

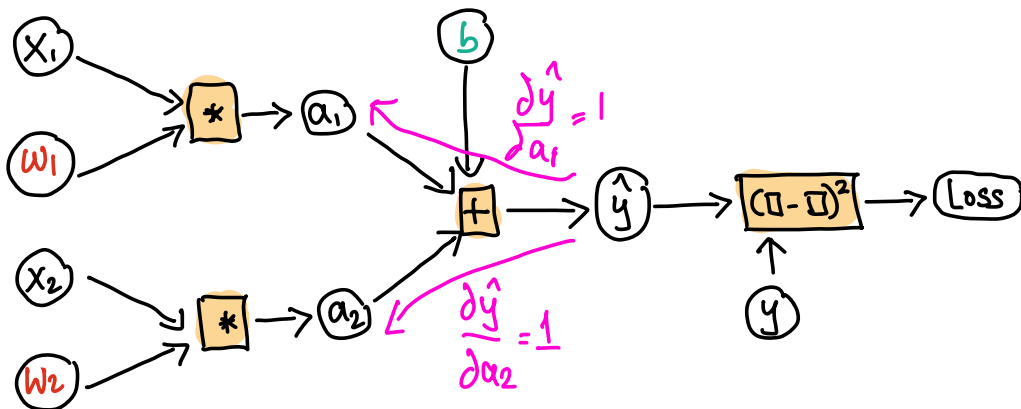
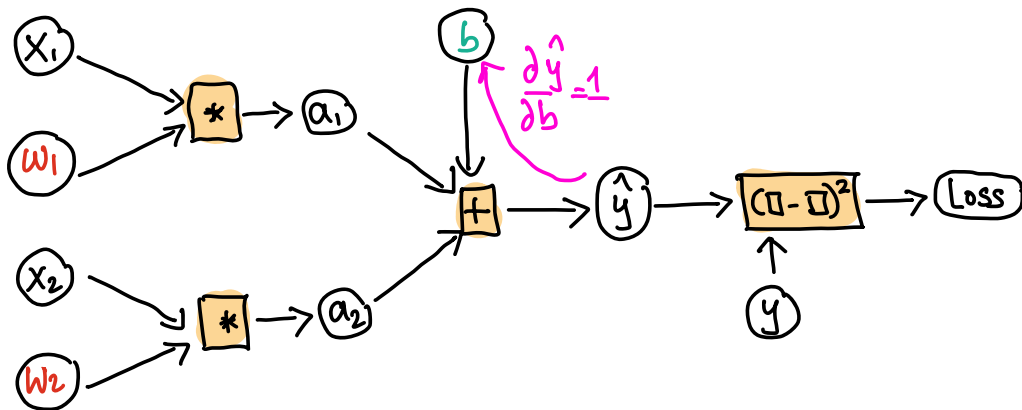
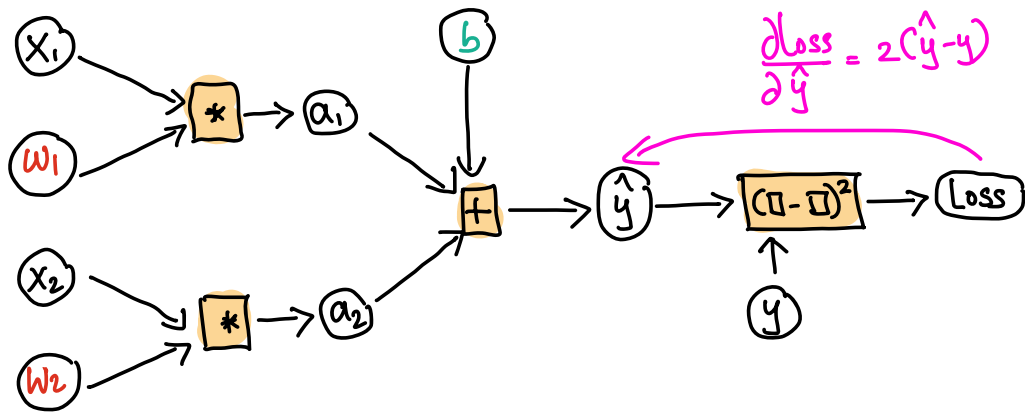
$$\begin{aligned}
 a_1 &= w_1 x_1 \\
 a_2 &= w_2 x_2 \\
 \hat{y} &= b + a_1 + a_2 \\
 \text{Loss} &= (\hat{y} - y)^2
 \end{aligned}$$

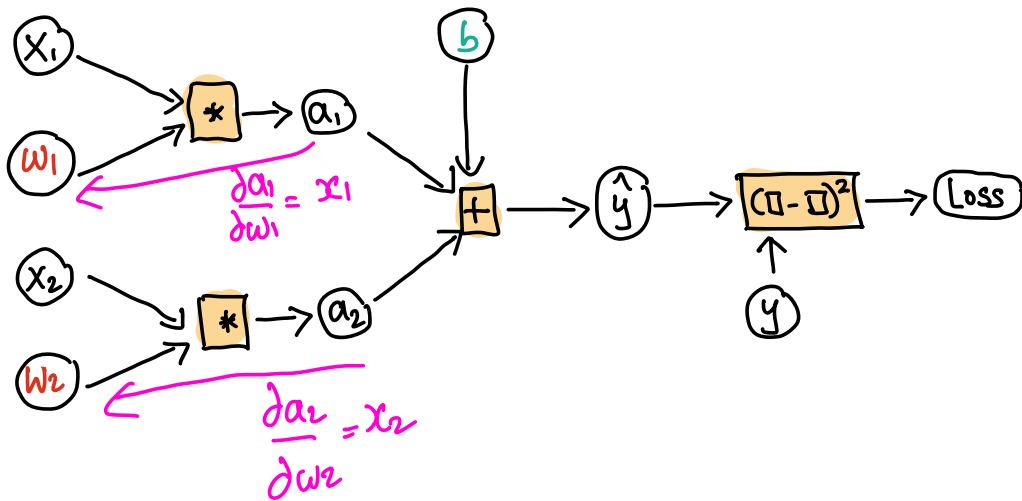


We will rewrite these equations as a **computational graph**

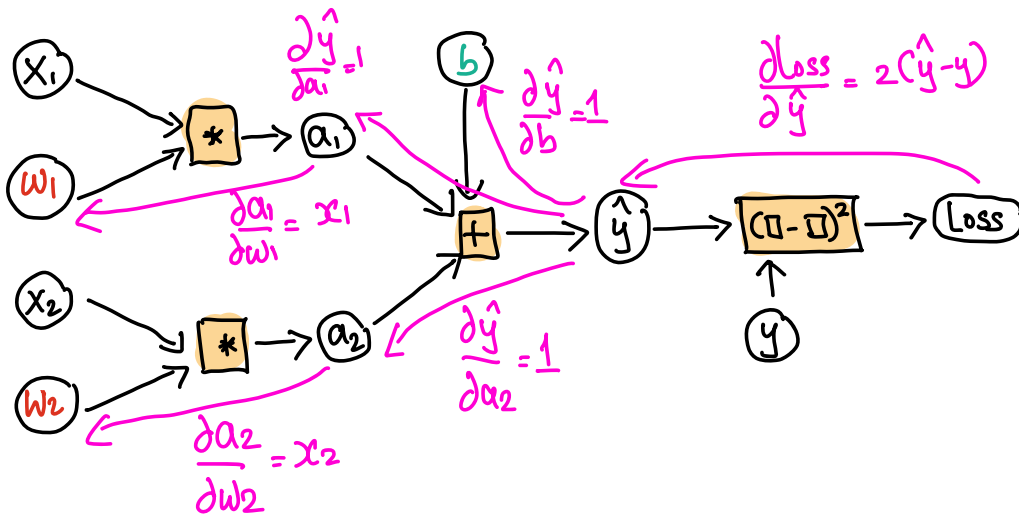


We can "attach" each of those little derivatives we calculated earlier to the graph.





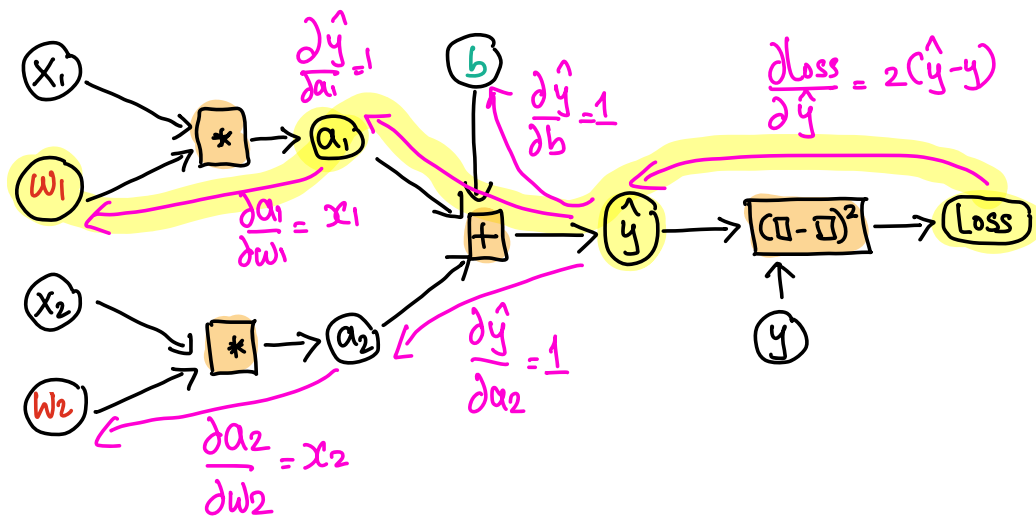
Putting everything together:



To calculate $\frac{\partial \text{Loss}}{\partial \text{any parameter}}$, start from the loss and travel backwards to the parameter, multiplying the partial derivatives as you go.

This is called **BACKPROPAGATION**

To calculate $\frac{\partial \text{loss}}{\partial w_1}$, for example:



Multiplying all the partial derivative on the yellow path, we get the answer:

$$\frac{\partial \text{loss}}{\partial w_1} = \frac{\partial \text{loss}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1} = 2(\hat{y} - y) \cdot 1 \cdot x_1$$

Does this match what we calculated earlier?

YES!!

Backprop is very efficient

- Calculate once and use many times (e.g. $\frac{\partial \text{loss}}{\partial g}$)
- (When more than one neuron in the layer) traversing backward is a series of matrix multiplications
- GPUs are perfect for matrix multiplications!!

