

Backpropagation via Computation Graphs

A Complete Worked Example for a 2-Layer Neural Network

1. Overview

This document walks through backpropagation for a simple 2-layer neural network (2 inputs, 2 hidden units, 1 output). We treat the network as a **computation graph** where values flow forward and gradients flow backward.

2. The Computation Graph

The network can be visualised as a chain of operations:

$$x \rightarrow z_1 \rightarrow a_1 \rightarrow z_2 \rightarrow a_2 \rightarrow L$$

(W1) (s) (W2) (s) (loss)

More explicitly, in vector form:

$$\begin{aligned} x \rightarrow z_1 &= W1 * x + b1 \rightarrow a_1 = \sigma(z_1) \\ &\quad \downarrow \\ z_2 &= W2 * a_1 + b2 \rightarrow a_2 = \sigma(z_2) \rightarrow L \end{aligned}$$

3. Forward Pass: Values Flow Left to Right

During the forward pass, we compute and **store** all intermediate values. These cached values will be needed during backpropagation.

Step-by-step computations:

Step	Node	Computation
1	Linear (hidden)	$z_1 = W1 * x + b1$
2	Activation (hidden)	$a_1 = \sigma(z_1)$
3	Linear (output)	$z_2 = W2 * a_1 + b2$
4	Activation (output)	$a_2 = \sigma(z_2)$
5	Loss	$L = (1/2) * (a_2 - y)^2$

4. Backward Pass: Gradients Flow Right to Left

Each node asks: "How does the loss change if I change slightly?" We apply the chain rule at each node, propagating gradients backward.

4.1 Loss to Output Activation

Starting at the loss node:

$$dL/da_2 = a_2 - y$$

4.2 Output Activation to Output Pre-activation

Since $a_2 = \sigma(z_2)$, by the chain rule:

$$da_2/dz_2 = \sigma'(z_2)$$

Combining these gives the **output error signal**:

$$\delta_2 = dL/dz_2 = (dL/da_2) * (da_2/dz_2) = (a_2 - y) * \sigma'(z_2)$$

4.3 Output Pre-activation to Weights, Bias, and Hidden Activations

Since $z_2 = W_2^T a_1 + b_2$, the gradient branches to three destinations:

Destination	Gradient Formula
Weight W_2	$dL/dW_2 = \delta_2 * a_1^T$
Bias b_2	$dL/db_2 = \delta_2$
Hidden activation a_1	$dL/da_1 = W_2^T * \delta_2$

4.4 Hidden Activation to Hidden Pre-activation

Since $a_1 = \sigma(z_1)$:

$$da_1/dz_1 = \sigma'(z_1)$$

Applying the chain rule node-by-node gives the **hidden error signal**:

$$\delta_1 = dL/dz_1 = (W_2^T * \delta_2) [elementwise*] \sigma'(z_1)$$

This is where δ_1 comes from: it is the gradient at the z_1 node, computed by propagating δ_2 backward through the weight matrix and element-wise multiplying by the local derivative of the activation.

4.5 Hidden Pre-activation to First-Layer Parameters

Since $z_1 = W_1^T x + b_1$:

Parameter	Gradient Formula
Weight W_1	$dL/dW_1 = \delta_1 * x^T$

Bias b1	$dL/db1 = \delta_1$
---------	---------------------

5. Summary of Error Signals

Symbol	Node	Formula	Meaning
δ_2	z_2	$(a_2 - y) * \sigma'(z_2)$	Output layer error
δ_1	z_1	$(W_2^T * \delta_2) [elem^*] \sigma'(z_1)$	Hidden layer error

6. Key Insights

δ_1 is not arbitrary

It is precisely dL/dz_1 - the gradient of the loss with respect to the hidden layer's pre-activation. It arises naturally from traversing the computation graph backward.

Backprop is reverse graph traversal

The forward pass computes values; the backward pass computes derivatives. Same graph, opposite direction.

Every delta corresponds to a node

Each error signal δ_k represents the gradient at a specific pre-activation node z_k . This is why error signals propagate - they follow the graph structure.

"Backpropagation is just the chain rule applied to a computation graph, one node at a time, from right to left."