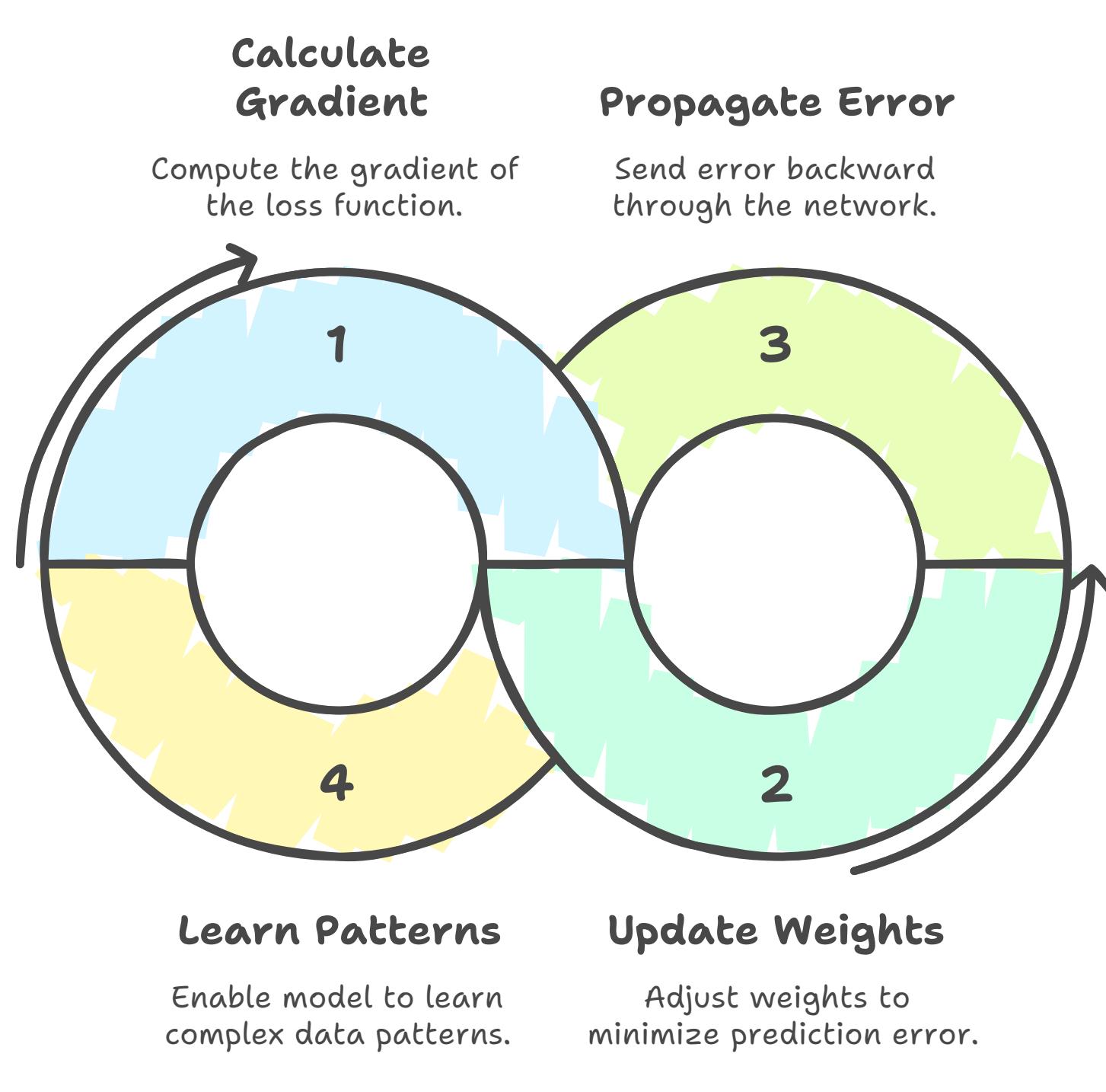


Backpropagation Cycle in Transformers



Transformer Decoder Process Flow

1. The decoder processes one token at a time. Let's use the case for the token "Le", and its decoder output is:

$$z = [1.07006, 0.95765]$$

3. Final decoder output: $z = [1.07006, 0.95765]$

$$\text{Vocabulary projection matrix: } W_{\text{vocab}} = \begin{bmatrix} 0.8 & 0.2 \\ 0.1 & 0.5 \\ -0.2 & 0.4 \end{bmatrix}$$

$$\text{logits} = z \times W_{\text{vocab}}^T = [1.0476, 0.5858, 0.1690]$$

$$\text{Softmax: } \hat{y} = [1.0476, 0.5858, 0.1690]$$

$$\text{Target: "chat" (index = 1), so:} \\ y = [0, 1, 0] \\ \text{Loss} = -\log(0.308) \approx 1.176$$

4. Linear Projection \rightarrow Softmax

$$\text{Step 1: Gradient of loss w.r.t. logits: } \partial \text{logits} / \partial L = y_{\text{hat}} - y = [0.489, -0.692, 0.203]$$

$$\text{Step 2: Gradient of loss w.r.t. } W_{\text{vocab}}: \partial L / \partial W_{\text{vocab}} = z^T \otimes \partial \text{logits} / \partial L \\ = \begin{bmatrix} 1.07006 \times 0.489, 1.07006 \times (-0.692), 1.07006 \times 0.203 \\ 0.95765 \times 0.489, 0.95765 \times (-0.692), 0.95765 \times 0.203 \end{bmatrix} = \begin{bmatrix} 0.523, -0.741, 0.217 \\ 0.468, -0.662, 0.194 \end{bmatrix}$$

$$\text{Step 3: Gradient w.r.t. decoder output } z: \\ \partial L / \partial z = \partial \text{logits} / \partial L \times W_{\text{vocab}}^T = [1.0476, 0.5858, 0.1690] \times \begin{bmatrix} 0.8 & 0.2 \\ 0.1 & 0.5 \\ -0.2 & 0.4 \end{bmatrix} \\ = [0.281, -0.167]$$

5. Decoder FFN

FFN uses two layers:

$$W1 \in \mathbb{R}^{(2 \times 4)} \\ W2 \in \mathbb{R}^{(4 \times 2)}$$

Given: Decoder input to FFN: $x = [0.5, 0.5]$

$$\text{From the forward pass: } z1 = x \otimes W1 = [0.4, 0.1, 0.35, -0.05] \\ a1 = \text{ReLU}(z1) = [0.4, 0.1, 0.35, 0] \\ z2 = a1 \otimes W2 = z \text{ (- passed to projection)}$$

$$\text{Step 1: Backprop through Linear 2: } \partial L / \partial W_2 = a1^T \otimes \partial L / \partial z \\ = \begin{bmatrix} 0.40 \times 0.281, 0.40 \times (-0.167) \\ 0.10 \times 0.281, 0.10 \times (-0.167) \\ 0.35 \times 0.281, 0.35 \times (-0.167) \\ 0.00 \times 0.281, 0.00 \times (-0.167) \end{bmatrix} = \begin{bmatrix} 0.1124, -0.0668 \\ 0.0281, -0.0167 \\ 0.0984, -0.0585 \\ 0.0000, 0.0000 \end{bmatrix}$$

$$\text{Step 2: Backprop to ReLU: } \partial L / \partial a1 = \partial L / \partial z \cdot W_2^T \\ = [0.281, -0.167] \cdot \begin{bmatrix} 0.4 & 0.1 \\ 0.1 & 0.35 \\ -0.2 & 0.4 \end{bmatrix} \\ = [0.146, 0.001, 0.207, -0.162]$$

Apply ReLU derivative (0 if $z1 < 0$): $\partial L / \partial z1 = [0.146, 0.001, 0.207, 0]$

$$\text{Step 3: Gradient of: } W1 = \partial L / \partial W1 = x^T \otimes \partial L / \partial z1 \\ \begin{bmatrix} 0.5 \times 0.146, 0.5 \times 0.001, 0.5 \times 0.207, 0 \\ 0.5 \times 0.146, 0.5 \times 0.001, 0.5 \times 0.207, 0 \end{bmatrix} = \begin{bmatrix} 0.073, 0.0005, 0.1035, 0 \\ 0.073, 0.0005, 0.1035, 0 \end{bmatrix}$$

6. Cross-Attention (Decoder \rightarrow Encoder)

$$\text{Decoder query: } Qd = [1.204, 0.895] \quad \text{Encoder keys: } K1 = [4.689, 1.905] \\ K2 = [4.818, 2.409]$$

$$\text{Softmax attention weights: } \alpha = [0.507, 0.493]$$

$$\text{Output of cross-attention: } o = \alpha1 \cdot V1 + \alpha2 \cdot V2 = [5.0918, 2.3458]$$

$$\text{Let gradient from FFN be: } \partial L / \partial o = [0.281, -0.167]$$

Step 1: Gradient w.r.t. values (V):

$$\partial L / \partial V1 = \alpha1 \otimes \partial L / \partial o \\ = [0.507] \cdot [0.281, -0.167] \\ = [0.1425, -0.0847] \\ \partial L / \partial V2 = \alpha2 \otimes \partial L / \partial o \\ = [0.493] \cdot [0.281, -0.167] \\ = [0.1385, -0.0823]$$

Step 2: Gradient w.r.t. attention weights:

$$\partial L / \partial \alpha1 = \partial L / \partial o \cdot V1 \\ = [0.281 \times 4.689 - 0.167 \times 1.905] \\ \approx 0.9994 \\ \partial L / \partial \alpha2 = \partial L / \partial o \cdot V2 \\ = [0.281 \times 4.818 - 0.167 \times 2.409] \\ \approx 0.9515$$

Use Softmax Jacobian for:

$$\frac{\partial \text{score} / \partial Q}{\partial \text{score} / \partial K} \\ \text{Then: } \partial L / \partial Q = \sum \partial L / \partial \alpha_i \cdot K_i \\ \partial L / \partial K_i = \sum \partial L / \partial \alpha_i \cdot Q_d$$

7. Decoder Self-Attention

From PDF:

$$\text{Single token: } x = [0.5, 0.5] \\ x = [0.5, 0.5]$$

$$Wq = I, Wk = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}, \\ Wv = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \\ \text{Output} = v \text{ (since attention score is 1.0)}$$

$$\text{Let } d_{\text{out}} = [dz1 \text{ from FFN}] \\ = [0.281, -0.167]$$

Then:

$$\partial L / \partial v = d_{\text{out}} \\ \partial L / \partial Wv = x^T \cdot d_{\text{out}} = [0.5, 0.5]^T \otimes [0.281, -0.167]$$

8. FINAL GRADIENT UPDATES

Assume learning rate $\eta = 0.1$:

$$\text{Projection: } W_{\text{vocab}} = W_{\text{vocab}} - 0.1 \times \partial W_{\text{vocab}}$$

$$\text{Decoder FFN: } W_1 = W_1 - 0.1 \times \partial W_1 \\ W_2 = W_2 - 0.1 \times \partial W_2$$

$$\text{Cross-Attention: } W_1 = W_1 - 0.1 \times \partial W_1 \\ W_2 = W_2 - 0.1 \times \partial W_2$$

Transformer Encoder Process Flow

1. Encoder Backpropagation

From decoder's cross-attention, we had:

$$\text{Query from decoder: } Q_d = [1.204, 0.895]$$

$$\text{Encoder key-value pairs: } K1 = [4.689, 1.905] \\ K2 = [4.818, 2.409]$$

$$\text{Attention weights: } \alpha1 = [0.507, 0.493]$$

$$\text{Output: } o = \alpha1 \cdot V1 + \alpha2 \cdot V2 = [5.0918, 2.3458]$$

$$\text{Gradient flowing into encoder values (from decoder):} \\ \partial L / \partial o = [0.281, -0.167]$$

Step 1: Gradient w.r.t. Encoder Values (V):

We already computed:

$$dL / dV1 = \alpha1 \cdot \partial L / \partial o = 0.507 \times [0.281, -0.167] = [0.1425, -0.0847]$$

$$dL / dV2 = \alpha2 \cdot \partial L / \partial o = 0.493 \times [0.281, -0.167] = [0.1385, -0.0823]$$

Step 2: Gradient w.r.t. Encoder Keys (K):

We decoder:

$$dL / dK1 = \partial L / \partial o \cdot V1 = 0.281 \times 4.689 - 0.167 \times 1.905 = 1.318$$

$$dL / dK2 = \partial L / \partial o \cdot V2 = 0.281 \times 4.818 - 0.167 \times 2.409 = 1.26$$

Assuming softmax with two scores $s1, s2$, we apply the Jacobian:

$$\text{Let } \alpha1 = 0.507, \alpha2 = 0.493$$

Then softmax gradients w.r.t. scores are:

$$dL / ds1 = d\alpha1 \cdot \alpha1(1-\alpha1) - d\alpha2 \cdot \alpha1\alpha2$$

Now backprop through:

$$s_i = Q_d \cdot K_i^T / \text{square_root}(d_k) \\ \Rightarrow dL / dK_i = (Q_d \cdot dL / ds_i) / \text{square_root}(d_k)$$

Step 3: Gradients w.r.t. Encoder Self-Attention:

$$\text{Each encoder token attends to others: } Q = X \cdot Wq, K = X \cdot Wk, V = X \cdot Wv$$

$$\text{Attention output for token}_i: \sum_j \text{softmax}(Q_{ik} \cdot V_j^T) V_j$$

We apply:

$$dL / dV = \alpha t \cdot d_{\text{out}} \quad \text{Then Compute: } dWq = X^T \cdot dQ \\ dL / dV = Q^T \cdot d_{\text{score}} \quad dWk = X^T \cdot dK \\ dL / dQ = d_{\text{score}} \cdot K \quad dWv = X^T \cdot dV$$

Step 4: Encoder Feedforward Network:

Assume input to encoder FFN is $x = [1, 1]$ for token 3:

Forward:

$$z1 = x \cdot W1 = [0.88, 0.22, 0.77, -0.11]$$

$$a1 = \text{ReLU}(z1) = [0.88, 0.22, 0.77, 0]$$

$$z2 = a1 \cdot W2 = \text{passed to next layer}$$

Backward:

$$dL / dz2 = [0.281, -0.167] \Rightarrow dW2 = a1^T \cdot dz2$$

$$da1 = dz2 \cdot W2^T, dz1 = da1 \cdot \text{ReLU}'(z1)$$

$$dW1 = x^T \cdot dz1$$

Step 5: Encoder LayerNorm + Residuals:

Let input to LayerNorm = x

$$\text{Mean: } \mu = \sum(x_i^2) / d$$

$$\text{Variance: } \sigma^2 = \sum((x_i - \mu)^2) / d$$

Backprop:

$$dx_{\text{norm}} = y \cdot \delta$$

$$\text{Use: } dx = [\delta - \text{mean}(\delta) - x_{\text{norm}} \cdot \text{mean}(\delta \cdot x_{\text{norm}})] / \sqrt{d + \epsilon}$$

Update LayerNorm weights:

$$dy = \sum \delta \cdot x_{\text{norm}}$$

$$d\beta = \sum \delta$$

Step 6: Final Encoder Weight Updates:

$$\text{With learning rate } \eta = 0.1: W_{\text{new}} = W_{\text{old}} - 0.1(dL / dW)$$