

Data matrix: $X = \begin{bmatrix} - & x_1 & - \\ & \vdots & \\ - & x_m & - \end{bmatrix}$; W_1 matrix = $\begin{bmatrix} - & w_1 & - \\ & \vdots & \\ - & w_{10} & - \end{bmatrix}$; W_2 matrix = $\begin{bmatrix} - & w_1^{(2)} & - \\ & \vdots & \\ - & w_{10}^{(2)} & - \end{bmatrix}$; $b_1 = \begin{bmatrix} b_1^1 \\ \vdots \\ b_{10}^1 \end{bmatrix}$; $b_2 = \begin{bmatrix} b_1^2 \\ \vdots \\ b_{10}^2 \end{bmatrix}$

Forward Pass

Input \rightarrow 1st layer

Weights: 754 weights
Samples: $x_1 \dots x_m$ features
Activations: $z_1 \dots z_m$ 10 nodes $10 \times m$

$$\begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix} = \begin{bmatrix} - & w_1 & - \\ & \vdots & \\ - & w_{10} & - \end{bmatrix} \cdot \begin{bmatrix} x_1 & \dots & x_m \end{bmatrix} + \begin{bmatrix} b_1 & \dots & b_m \end{bmatrix}$$

$10 \times m = 10 \times 754 + 10 \times m$

Focus on how one sample contributes to all 10 nodes

$$\begin{bmatrix} - & w_1 & - \\ & \vdots & \\ - & w_{10} & - \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_1 \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ n_{10} \end{bmatrix}$$

$10 \times 754 \quad 754 \times 1 = 10 \times 1$

Same row vector of biases m times

Activations: $z_1 \dots z_m$ 10 nodes $10 \times m$

Each sample produces its own set of 10 values corresponding to the 10 nodes in layer one.

\rightarrow obviously that's the case. You only combine the collective error of all the samples when computing loss. But for prediction obviously there will be m unique predictions corresponding to the m samples.

Each prediction is a list of 10 probabilities indicating likelihood of each possible answer.

$$A^{(1)}_{10 \times m} = \text{ReLU} \left(\begin{bmatrix} \text{activations} \end{bmatrix} \right)$$

\rightarrow Apply ReLU to every node for every sample

Math Form:

$$z' = w' x^T + b$$

$10 \times m \quad 10 \times m \quad 754 \times m \quad (10 \times 1) \rightarrow 10 \times m$

$$A'_{10 \times m} = \text{ReLU}(z') ; A_{k,i}^{(1)} = \max(0, z_{k,i}^{(1)})$$

(each individual neuron goes through ReLU to determine whether it will contribute to the next layer. * This idea becomes important again during backprop. *)

2nd Layer:

Treat z' from previous layer as input for second layer. It works out nicely because the data already arranged

$$z^2 = w^2 A' + b^2$$

$10 \times m \quad 10 \times 10 \quad 10 \times m \quad (10 \times 1) \rightarrow 10 \times m$

Same 10×1 column m times

$$\begin{bmatrix} z_1^{(2)} & \dots & z_m^{(2)} \end{bmatrix} = \begin{bmatrix} - & w_1^{(2)} & - \\ & \vdots & \\ - & w_{10}^{(2)} & - \end{bmatrix} \begin{bmatrix} A_1^{(1)} & \dots & A_m^{(1)} \end{bmatrix} + \begin{bmatrix} b^{(2)} & \dots & b^{(2)} \end{bmatrix}$$

$10 \times m \quad 10 \times 10 \quad 10 \times m \quad 10 \times m$

Final output

$$A^{(2)} = \text{softmax} \left(\begin{bmatrix} \text{activations}^{(2)} \end{bmatrix} \right)$$

$A^{(2)}$ in the form of probabilities

\rightarrow Pick largest one to make decision

Backpropagation

1. Compute loss function

one hot encode labels

ex) $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ ← correct answer

Calculate: sparse-categorical cross entropy loss (for one sample)

↳ The labels will not be one hot - encoded

$$\text{Loss}^{(i)} = -\log(\hat{y}_{y^{(i)}}^{(i)})$$

↳ The predicted softmax prob. for the correct label.

• \hat{y} = predicted softmax vector (A^2 in forward pass notation)

$$\text{Cost} = -\frac{1}{m} \sum_{i=1}^m \log(\hat{y}_{y^{(i)}}^{(i)})$$

Cost in my notation:

$$\text{cost} = -\frac{1}{m} \sum_{i=1}^m \log(A_{y^{(i)}}^{2(i)})$$

← prediction for sample i corresponds to i th column of A^2 .

$$\text{loss}^i = -\log(A_{y^{(i)}}^{2(i)}) ; \text{loss computed for } i\text{th sample}$$

Gradient:

*** The math for gradient step done as though you had one-hot labels. But in code

no need to form full $10 \times m$ Y matrix of one-hot-encoded vector for each label.

Instead for each sample you just subtract 1 at the true-class labels index. See, no need to store a $10 \times m$ Y matrix.

$$\frac{\partial C}{\partial Z^2} = \frac{1}{m} (A^2 - Y)$$

↳ in this expression this is a $10 \times m$ matrix with each col being one-hot encoded label matrix for a given sample. Code implementation will not use this!

Update W_2 :

$$\frac{\partial C}{\partial W^2} = \frac{\partial C}{\partial Z^2} \cdot \frac{\partial Z^2}{\partial W^2} = \frac{1}{m} (A^2 - Y) \cdot (A^1)$$

$$\text{Grad descent: } \underset{10 \times 10}{W^{2, \text{new}}} = \underset{10 \times 10}{W^{2, \text{old}}} - \eta \nabla C$$

$$\underset{10 \times m}{\frac{1}{m} (A^2 - Y)} \cdot \underset{10 \times m}{(A^1)} \approx \underset{10 \times 10}{W^2} \rightarrow \boxed{\underset{10 \times m}{\frac{1}{m} (A^2 - Y)} \underset{m \times 10}{(A^1)^T}} = \underset{10 \times 10}{W^2}$$

Update b^2 :

$$\frac{\partial C}{\partial b^2} = \frac{\partial C}{\partial Z^2} \cdot \frac{\partial Z^2}{\partial b^2} = \frac{1}{m} (A^2 - Y) (1)$$

Recall the b^2 matrix was formed by repeating the same b vector m times. We need to aggregate the contribution to gradient from every sample into a single (10×1) vector to match dimension of b^2 which is really just a 10×1 vector. Contribution of each sample to grad w respect to b_2 is stored as col of $\frac{1}{m} (A^2 - Y)$ computed above.

$$\frac{\partial C}{\partial b^2} = \frac{1}{m} \sum_{i=1}^m (A_{:,i}^2 - Y_{:,i})$$

↳ means add up the cols of $(A^2 - Y)$ to get from $10 \times m$ to 10×1 vector.

Update w_1 :

$$\frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial z^2} \cdot \frac{\partial z^2}{\partial A'} \cdot \frac{\partial A'}{\partial z^1} \cdot \frac{\partial z^1}{\partial w_1}$$

1. $dZ^{(2)} = \frac{1}{m} (A^{(2)} - Y)$ ex) $dz^2 = \frac{\partial C}{\partial z^2}$
2. $dA^{(1)} = W_2^T dZ^{(2)}$
3. $dZ^{(1)} = dA^{(1)} \odot (Z^{(1)} > 0)$ (element-wise ReLU mask)
4. $dW_1 = \frac{\partial C}{\partial w_1} dz^{(1)} X$

ReLU:

$$A'_{k,i} = \max(0, z_{k,i}) \quad \text{turns off neurons that won't contribute to next layer.}$$

ReLU mask:

Why apply the mask and what does it actually do?

Why:

When you back-propagate through a nonlinearity like ReLU, you need to ask "for each neuron, did it actually contribute to the forward pass?" The ReLU mask is how you answer that.

Derivative of ReLU:

$$\frac{d}{dz_{k,i}} \max(0, z_{k,i}) = \begin{cases} 1, & z_{k,i} > 0 \\ 0, & z_{k,i} \leq 0 \end{cases}$$

↳ This restricts gradient flow to the neurons that contributed to the second layer during forward pass.

The mask:

$(z^{(1)} > 0)$ is a $10 \times m$ matrix of 0's and 1's indicating active or off neurons.
A bit like A' but instead of having $z_{k,i}$ where $z_{k,i} > 0$ we have a 1.

What does the mask do?

You've already got the "raw" gradient into that layer, $dA^{(1)}$. To turn it into the gradient w.r.t. the pre-activation $Z^{(1)}$, you do an elementwise multiply:

$$dZ^{(1)} = dA^{(1)} \odot (Z^{(1)} > 0).$$

- Wherever the mask is 1, $dZ^{(1)}_{k,i} = dA^{(1)}_{k,i}$ (gradient flows unchanged). ← means gradient flows back to weights of neurons that contributed to the forward propagation only.
- Wherever the mask is 0, $dZ^{(1)}_{k,i} = 0$ (no gradient flows back into that neuron).

$$\frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial z^2} \cdot \frac{\partial z^2}{\partial A'} \cdot \frac{\partial A'}{\partial z^1} \cdot \frac{\partial z^1}{\partial w_1} = \frac{1}{m} \left[\underset{10 \times 10}{(W_2^2)^T} \underset{10 \times m}{(A^{(2)} - Y)} \odot \underset{10 \times m}{(z^{(1)} > 0)} \right] \underset{m \times 784}{(X)}$$

Note: w_1 is 10×784

1. $dZ^{(2)} = \frac{1}{m} (A^{(2)} - Y)$ ex) $dz^2 = \frac{\partial C}{\partial z^2}$
2. $dA^{(1)} = W_2^T dZ^{(2)}$
3. $dZ^{(1)} = dA^{(1)} \odot (Z^{(1)} > 0)$ (element-wise ReLU mask)
4. $dW_1 = \frac{\partial C}{\partial w_1} dz^{(1)} X$

Update b_1 :

$$\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial z^2} \cdot \frac{\partial z^2}{\partial A'} \cdot \frac{\partial A'}{\partial z^1} \cdot \frac{\partial z^1}{\partial b_1}$$

$$\frac{\partial C}{\partial b_1} = \sum_{i=1}^m \frac{1}{m} \left[\underset{10 \times 10}{(W_2^2)^T} \underset{10 \times m}{(A^{(2)} - Y)} \odot \underset{10 \times m}{(z^{(1)} > 0)} \right] \quad \left(\text{sum up each column of this matrix leaving you with a } 10 \times 1 \text{ matrix.} \right)$$

Additional Notes:

1. Be attentive to the notation used for data matrix. Samson Zhang video uses X where cols are samples not rows.
2. Pay attention to summation VS matrix notation.