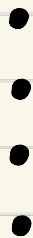
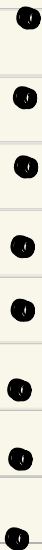


#1 Initialize the neural network

so we have → 9 input features
4 classification stages

for shallow neural network lets
consider 8 hidden units



input = 9

hidden = 8

output = 4



weight matrix (W) = $[9 \times 8]$

weight matrix (W) = $[8 \times 4]$

bias (b) = $[1 \times 8]$

bias (b) = $[1 \times 4]$

- Understanding matrix multiplication (input to hidden layer)

$$\text{weight matrix } (W) = [9 \times 8]$$

$$\text{input features } (x) = [620 \times 9]$$

number of samples (n) → 620
number of unique features → 9

$$\text{so } F(x) = X \cdot W + b$$

shape → $620 \times 8 + b[1 \times 8]$

the bias term will be broadcasted to match the shape

$$F(x) = (620, 8) \rightarrow \text{output shape of input to hidden layer}$$

• Understanding matrix multiplication (hidden to output layer)

$$F(x) = (620, 8) \rightarrow \text{output shape of input to hidden layer (calculated previously)}$$

$\hookrightarrow Z_{\text{hidden}}$

$$\text{weight matrix } (W) = [8 \times 4]$$

$$\text{bias } (b) = [1 \times 4]$$

$$F(x) = Z_{\text{hidden}} \cdot W + b$$

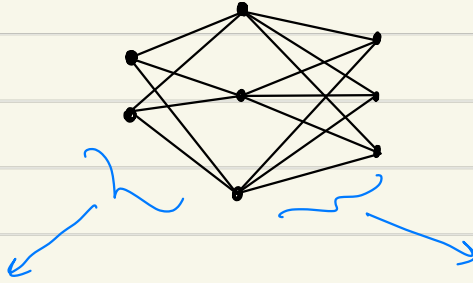
shape $(620, 8) \times (8 \times 4)$

$$620 \times 4 + b[1 \times 4]$$

$$F(x) = (620, 4)$$

\rightarrow = number of classification

#2 Forward Propagation



Input to Hidden layer

calculate

$$Z = X \cdot W + b$$

↓ pass it to activation function

$$A = \text{ReLU}(0, Z)$$

↳ $\max(0, Z)$

Hidden to Output layer

using activation A from previous layer

$$Z = A \cdot W + b$$

apply softmax function to convert logits to probabilities

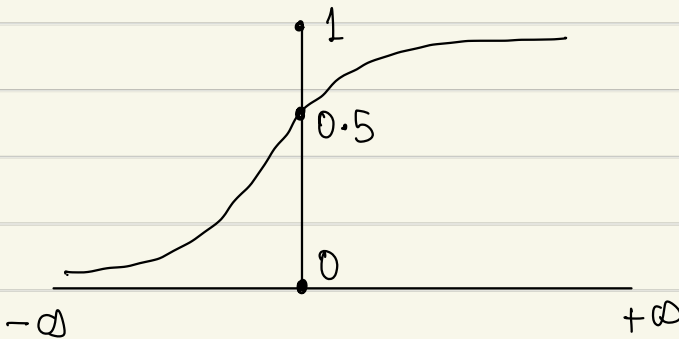
↳ basically converting numbers between the range of 0 and 1

- More about Softmax

Softmax is an activation function applied to the last layer of neural network

$$\sigma(z) = \frac{e^{z_i}}{\sum_i^k e^{z_i}}$$

k = number of classes



- Cross Entropy Loss Function

Now our network predicted classes for our input data
What Next?

check how correct is our neural network

\therefore we use Cross Entropy Loss

$$H = - \sum_{i=1}^n \underbrace{p(x)}_{\substack{\text{the correct} \\ \text{prediction} \\ \text{in our} \\ \text{dataset}}} \underbrace{\log p(x)}_{\substack{\text{prediction} \\ \text{our} \\ \text{network} \\ \text{made}}}$$

The diagram shows the Cross Entropy Loss function $H = - \sum_{i=1}^n p(x) \log p(x)$ enclosed in a box. Annotations include: an arrow from 'number of classes' to the summation index n ; an arrow from 'the correct prediction in our dataset' to the term $p(x)$; and an arrow from 'prediction our network made' to the term $\log p(x)$.

3. Backpropagation

Previously we calculated errors using Cross Entropy Loss. Now we need to backpropagate these errors to update weights and biases

- what is our goal in backpropagation?
 - calculate gradients of the loss wrt
 - 1) weights in the hidden and output layers
 - 2) biases in the hidden and output layers
 - take these gradients and update the weights and biases

1) output layer gradients

in here we calculate Δ_{output}

$$= (y_{\text{pred}} - y_{\text{true}})$$

→ shows us the direction in which correction is needed

2) hidden layer gradients

take output layer's error (Δ_{output}) and run it back through the network to understand each hidden neuron's contribution to error

in here we calculate Δ_{hidden}

$$\Delta_{\text{hidden}} = (\Delta_{\text{output}} \cdot W_{\text{hidden-output}}) \times (I_{\text{hidden}} > 0)$$

only active neurons
contribute to error

→ \therefore checks if ReLU activation is 0

3) calculate weight and bias gradients

a) Gradient for weight and bias updates in hidden to output layer

$$\nabla W_{\text{hidden-to-output}} = \frac{a_{\text{hidden}}}{\text{activations in hidden layer}} \cdot \frac{\Delta_{\text{output}}}{\text{output error}}$$

$$\nabla b_{\text{hidden-to-output}} = \sum \Delta_{\text{output}}$$

b) Gradient for weight and bias updates in hidden to input layer

$$\nabla W_{\text{input-to-hidden}} = \frac{X}{\text{input features}} \cdot \frac{\Delta_{\text{hidden}}}{\text{hidden neurons error}}$$

$$\nabla b_{\text{hidden-to-output}} = \sum \Delta_{\text{hidden}}$$

4) We use these gradients to update our trainable parameters

a) hidden to output layer

$$\text{new } W = \text{old } W - (\text{learning rate}) \times \nabla W_{\text{hidden-to-output}}$$

$$\text{new } b = \text{old } b - (\text{learning rate}) \times \nabla b_{\text{hidden-to-output}}$$

b) input to hidden layer

$$\text{new } W = \text{old } W - (\text{learning rate}) \times \nabla W_{\text{input-to-hidden}}$$

$$\text{new } b = \text{old } b - (\text{learning rate}) \times \nabla b_{\text{hidden-to-output}}$$

4. Stochastic Gradient Descent

1) Data shuffling

shuffle the dataset but why?

- prevent model from learning order-specific patterns
- now each batch represents the entire dataset randomly
- mainly reduces bias and improves generalization

2) Mini batch processing

- breaks large datasets into smaller and manageable chunks

3) Forward Propagation

→ input data flows through the network and it predicts labels

4) Loss Calculation

→ measure how far are the predictions from true labels

5) Compute Gradients

→ calculate gradients with the help of loss function for each trainable parameter

6) Update Parameters

→ use the calculated gradients and learning rate to update weights and biases