

BB1B - Neural Networks

(1)

- ① what is a neural network?
- ② Gradient Descent, how neural networks learn
- ③ Backpropagation, step-by-step
- ④ Backpropagation calculus

What is a neural network?

This video handles the structure of a NN - no learning

Plain NN = Multilayer perceptron

neuron = thing that holds a number

with image neuron = Pixel = activation

input \rightarrow hidden \rightarrow output

input num of neurons = Pixel count

hidden neuron = arbitrary choice

with Digit ~~recognition~~ recognition the output neurons correspond to a possible digit 0-9

think about neurons in one layer being active & setting of ~~of~~ ^{neurons} in the next layer ~~at~~

Connectors between neurons are weights

Activation would be better if between 0-1 (sigmoid)

(2)

Sigmoid is applied to weighted sum of neurons & weights

$$\sigma(w_1 a_1 + w_2 a_2 + \dots + w_n a_n)$$

bias is used to adjust activation of weighted sum ; i.e. if weight = 10 - 10 bias = 0 so neuron isn't active

Can think of each layer as learning new characteristic of the problem

Linear Algebra way to write weighted sums of ~~weights~~

neurons of input layer ($a^{(0)}$)

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$$

bias

Each row is the weights
for a connecting neuron

$$a^{(1)} \sigma(W a^{(0)} + b)$$

Much easier way to write the connections between layers

Python: $a = \text{Sigmoid}(\text{np.dot}(w, a) + b)$

A better way to view a neuron might be as a function rather than a number

(3)

Sigmoid is almost never used now

instead the main activation used is $\text{ReLU}(a)$

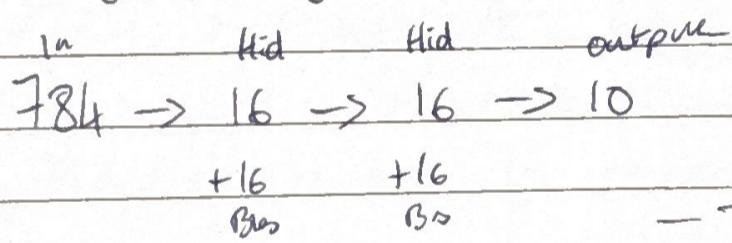
$$\text{ReLU} = \max(0, a)$$

i.e. if less than 0 then neuron is "turned off"

Gradient Descent, how NNs Learn?

- Grad Desc
- Analyze network

Digit recognition - 28×28 image = 784



Bk weights & bias
— These are the things to change

UNIST database = collection of hand written Digits

to start the training process initialize the network by randomly assigning all weights & bias

Outputs will be rubbish but use these output to create a loss function

$$0: (\frac{\text{opt}}{\text{pred} - \text{actual}})^2$$

in the case of 0-9 Digits 1:

2:

...

then do average cost over all training images

sum

(4)

This total or average cost function is the performance of the network that needs to be optimize against

NN function:

Input: 784 number (pixels)

Output: 10 number (0-9)

Params: 13,002 weights & biases

Cost func = Sums but 1 output metric
overall train data

Cost function = Performance

Now need a way to set w & b to improve perform

$C(w)$ = Cost & we create a start point
then we want to work out which way to
step in order to make $C(w)$ lower

How to change the input to the cost function: w & b

- find the slope of the function
- if positive slope, move left: y neg & move right
iterate to achieve minimum

global
Could land in any ~~local~~ minima

Gradient of a function gives ~~the~~ the direction
of steepest descent $\nabla C(u, y)$

(5)

the basic "algo" for ~~computing~~^{minimizing} the function:

- ① Compute gradient direction ∇C
- ② Small step in $-\nabla C$ direction
- ③ Repeat

the concept is the same for $-\nabla C(\vec{w})$

$$\vec{w} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$-\nabla C(\vec{w}) = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

grad of cost
w/ respect to weight

minus to reverse slope

w_n gets updated by $-\nabla C(w) = y_n$

Recall, the cost function is an average over all training data images

so an improvement to the cost is an improvement to the performance of the whole prediction

the algo for computing this gradient exactly is the main part of how NNs work

- Backprop
(Next video)

learn = min cost function

All of the weights to a neuron can be visualized to see the "pattern" of learning

Multi-layer perceptron is old tech from the 80s - 90s

But it is the precursor to more modern tech

Visualizing learning weights in NNs shows you just how basic & often not good the learning is

BackProp, Step-by-Step

the negative gradient of the cost function is what is used to update the weights & biases

$$-\nabla C(\dots)$$

Computing the gradient for the entire training set is very costly & slow so it is done in "mini-batches"

Update weights @ each batch (?)

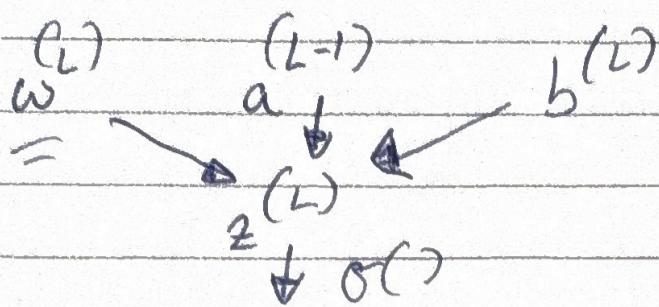
- Stochastic gradient descent

Backprop-Calculus

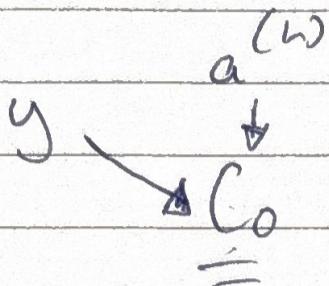
$$\text{a}^{(L-1)} \xrightarrow{\omega^{(L)}} \text{a}^{(L)} \quad \text{vs} \quad \begin{matrix} \text{a}^{(L)} \\ \text{y} \end{matrix} \quad C_0(\cdot) = (\text{a}^{(L)} - \text{y})^2$$

$$\text{a}^{(L)} = \omega^{(L)} \text{a}^{(L-1)} + b^{(L)} \quad \text{vs} \quad z^{(L)} \leftarrow \text{a}^{(L)} = \sigma(z^{(L)})$$

(7)



Connections of what is needed to calculate what



We want to know what happens to the cost func if we change a weight $w^{(L)}$ by a tiny amount

$$\frac{\partial C_0}{\partial w^{(L)}} \text{ w/ respect to } w^{(L)}$$

A change in $w^{(L)}$ changes $z^{(L)}$ changes $a^{(L)}$ changes C_0

there is a ^{chain}~~change~~ of intermediate changes that need to be follow to observe the change

$$\frac{\partial C_0}{\partial w^L} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \cdot \frac{\partial a^{(L)}}{\partial z_L} \frac{\partial C_0}{\partial a^{(L)}}$$

this is the chain rule - multiply these 3 ratios to understand how the cost change w/ respect to small changes in $w^{(L)}$



(8)

Working through the derivative chain:

$$* C_0 = (a^{(L)} - y)^2$$

$$* z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$* a^{(L)} = \sigma(z^{(L)})$$

$$\frac{\partial C_0}{\partial a^L} = 2(a^{(L)} - y)$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$\frac{\partial C_0}{\partial w^{(L)}} = a^{(L-1)} \cdot \sigma'(z^{(L)}) \cdot 2(a^{(L)} - y)$$

Note this is just single example in net = $\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum \frac{\partial C_k}{\partial w^{(L)}}$

* overall training examples

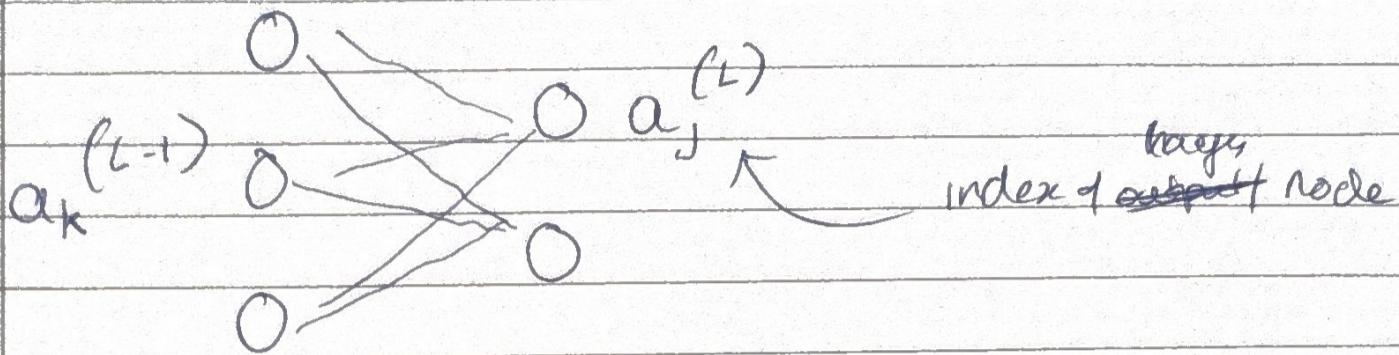
∇C = the true gradient is built up over all train examples & all partial derivatives of all weight and biases

The route is just the same for bias

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial z^L}{\partial b^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L}$$

(9)

the concept remains much the same as we increase the complexity of the network



$$C_0 = \sum (a_j^{(L)} - y_j)^2$$

Sum of nodes output errors across j tries

$w_{jk}^{(L)}$ = weights are denoted by their end, Stark nodes

$$\frac{\partial C_0}{\partial w_{jk}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}$$

Cost function for each weight in the network

This gives you the value through which to "update" the layer of ~~weig~~ weight

Note network gets updated sequentially

nearest \rightarrow Repest layers

recalc is required after each update starting from the cost function

weight change activation values which change cost func