

youtube

(1)

first principles of CV - Neural Networks Series

- (1) overview
- (2) Perceptron → Perceptron network multilayer P
- (3) Activation Function
- (4) Neural network
- (5) Gradient Descent
- (6) Backpropagation
- (7) Example applications

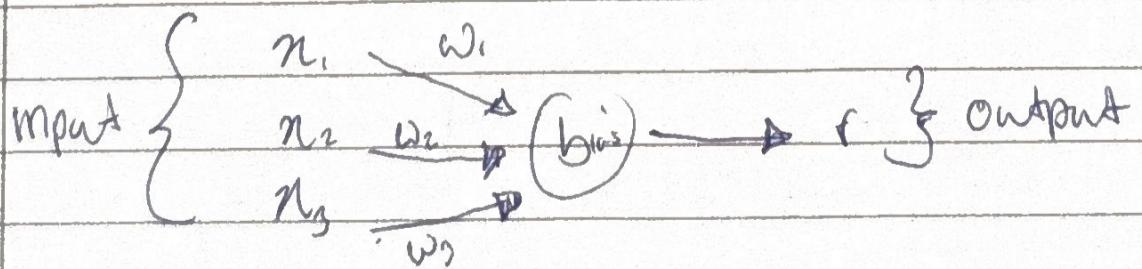
Overview

NN

- (1) Perceptron → very similar to linear model
- (2) Perceptron Network
- (3) Activation functions
- (4) Neural network → very big, how to train? GD
- (5) Backprop → innov because GD is xpen
Grad Desc is very comp
- (6) Example applications

Perceptron

very simple type of neuron



$$f = \begin{cases} 0 & \text{if } w \cdot n + b \leq 0 \\ 1 & \text{if } w \cdot n + b > 0 \end{cases}$$

2)

$$\text{let } z = w \cdot x + b$$

→ this is what activation func gets applied to

$$a = f(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$$

The activation function of a perceptron is a step function $(0, 1)$

Example Q: Will you go to movies

21. Is the weather good?

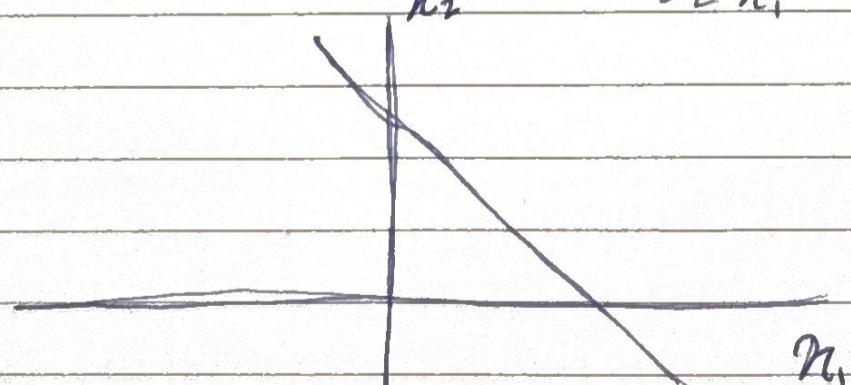
2 Do you have company?

n_3 is the theater nearby.

- Set weights $w_1 = 4$ $w_2 = 2$ $w_3 = 2$
 - Set bias = min threshold $b = -5$
i.e. if $w_1 \geq (w_2 \text{ or } w_3)$ active term ≥ 0

$$x_1 \xrightarrow{-2} \textcircled{3} \xrightarrow{} a = f(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$$

$$x_2 = -2x_1 - 2x_3 + 3$$



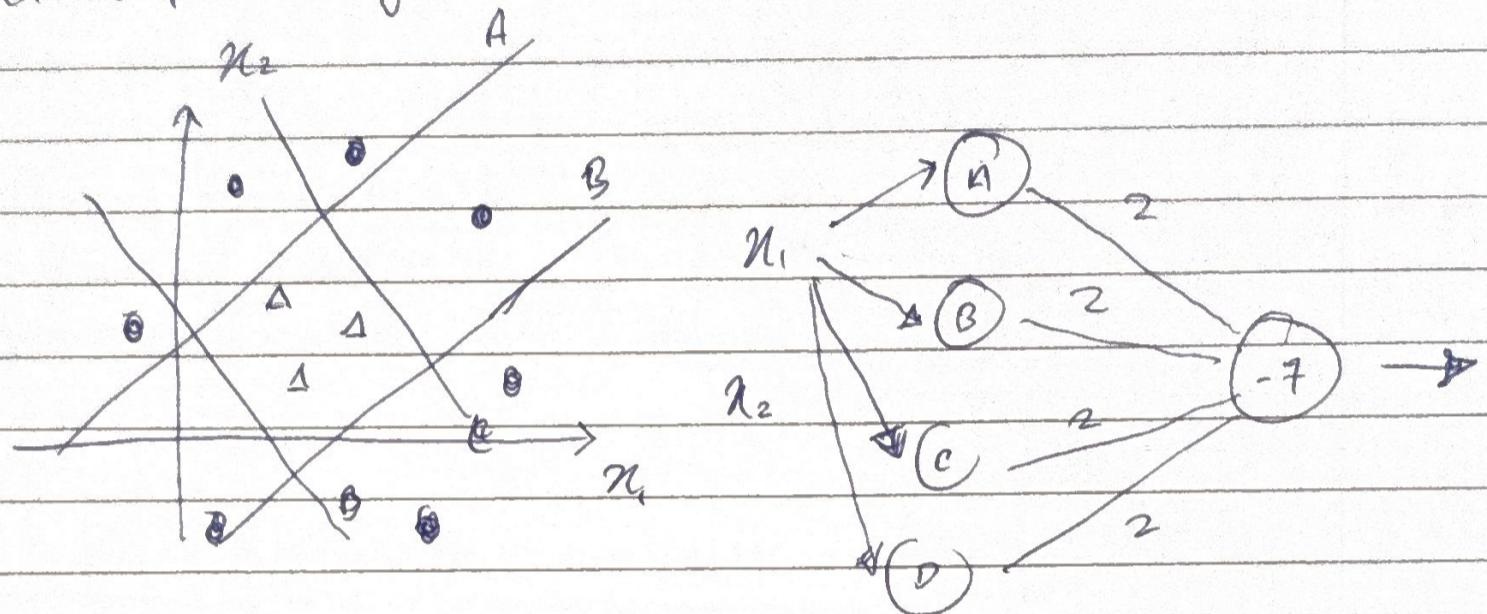
this is what the
inside of a person
is doing

A lin classifier

Perceptron Network

what can we do with a network of perceptrons?

A network allows us to construct a linear classifier for complex regions



- each node creates its own plane (ϕ_i)
- only when each node is switched on will the total cross the threshold

Multi-layer Perceptron Network

input layer \rightarrow hidden layers \rightarrow output layer

Activation function

- function that takes you from input \rightarrow output
- for perceptron = step function by default
- the task w/ perceptron is to find the correct weights & bias (parameters)
- How does a change in w or b change the output?
- this is an issue for the step function because of 0 or 1 output
- can see the impacts of changing w & b , only when 0, 1 change
- too unstable

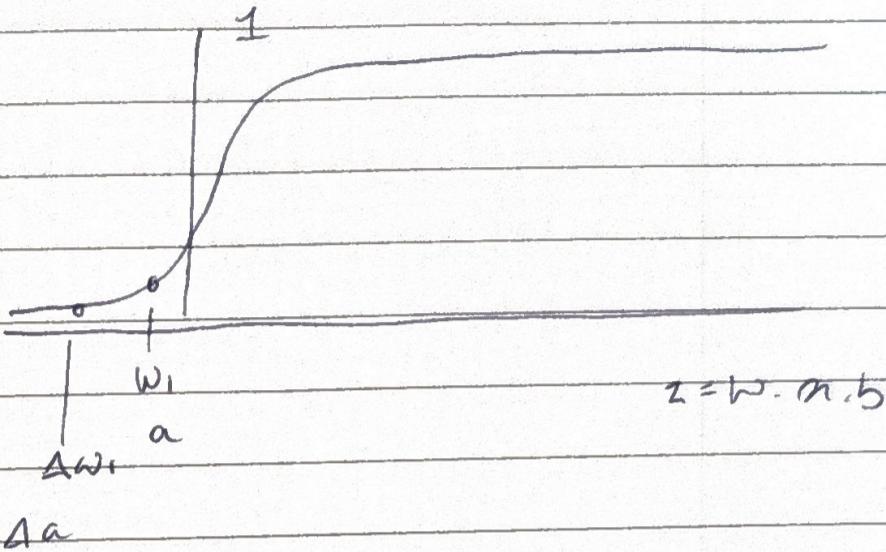
(4)

Step function is too dramatic, would like it to be
 A smoother function

Sigmoid Neuron

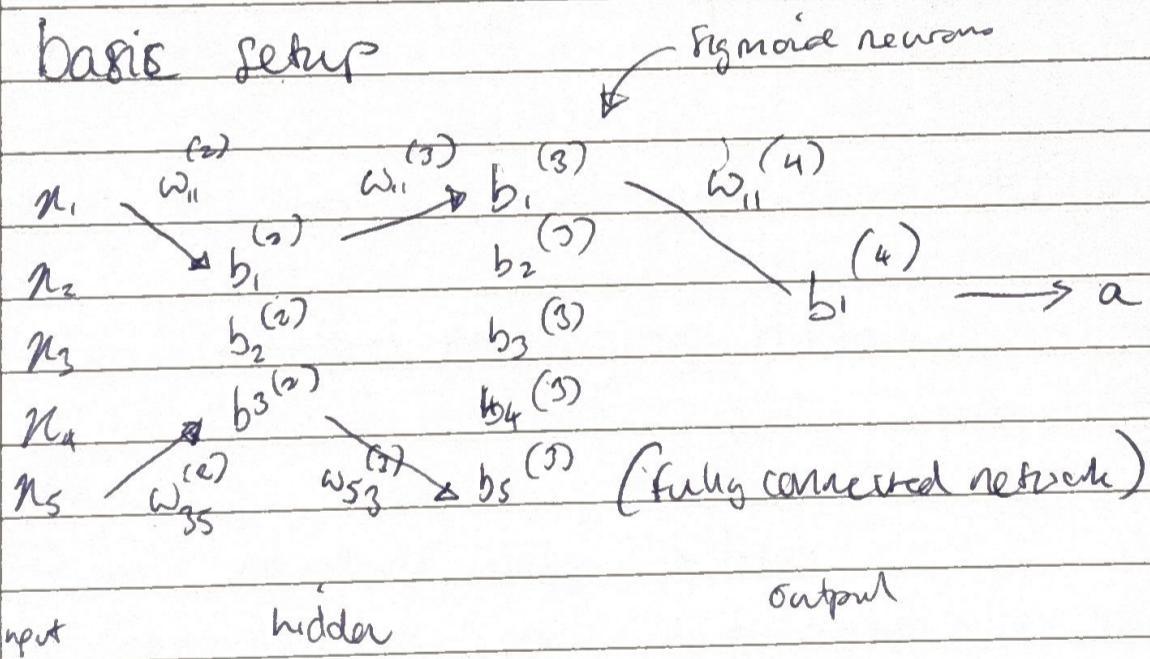
$$a = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Allows for measurable output w.r.t respect to change in w or b



Neural Networks - Construct & train

basis setup



Problem \rightarrow recognise handwritten

MNIST \rightarrow binarised, thresholded, bound 20x20 pixels

Output: Digits 0-9 (10)

within 28x28 image

goal = high activation output for digits

neilson 2015 has model for this

(3)

784

neuron

Neilsen used inputs of 784 neurons = a pixel for each pixel in the 28×28

which go into 30 hidden layer neurons

out into 10 output neurons for each digit

$$784 \rightarrow 30 \rightarrow 10$$

gets 95% Accuracy

if we don't have a trained ~~net~~ network we start by randomizing the weights & biases

this will ~~be~~ give a poor output & unpredictable outputs

how to change weights & bias to improve?

need labelled training data

Steps

- ① init random weights & bias
- ② Compute network activations (outputs) for each of the training images
- ③ Compute the cost for whole train data set

$$C_x(w, b) = \|\hat{a}(x) - a(x|w, b)\|^2$$

euclidean distances

how this looks
per image

actual	Pred	0	0.2
0	0	0.3	0.1
0	1	-0.0	0.9
0	0	0.1	0.2
0	0	0.2	0.0

$$\text{for whole data: } C(w, b) = \frac{1}{n} \sum C_x(w, b)$$

- ④ update weights of w & b to min cost

Gradient Descent

Adjust weights & biases using optimization process (GD)

Want to min cost function using paras (w, b)

Start w/ w & b values

- adjust until min cost

How do you measure change in a function?

Taylor Series

- This tells us \rightarrow

$$f(x + \Delta x, y + \Delta y) \approx f(x, y) + \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y$$

$$\Delta f(x, y) \approx \frac{\partial f}{\partial x} \Delta x + \frac{\partial f}{\partial y} \Delta y$$

Ultimately shows us a way to measure the change in a function w/ respect to changes in its paras

Applied to cost function:

$$\Delta C \approx \frac{\partial C}{\partial w_i} \Delta w_i + \frac{\partial C}{\partial b_j} \Delta b_j$$

$$\text{Vector form} = \Delta C = \begin{bmatrix} \frac{\partial C}{\partial w_1} & \frac{\partial C}{\partial b_1} \\ \frac{\partial C}{\partial w_2} & \frac{\partial C}{\partial b_2} \end{bmatrix} \cdot \begin{bmatrix} \Delta w_1 \\ \Delta b_1 \end{bmatrix}$$

∇C : gradient of C

ΔV : change in vars

Gradients /^{1st} derivatives

∇C

Para change

ΔV

$$\nabla C \cdot \Delta V = \Delta \text{Cost}$$

∇C gives you value to increase cost function, we want opposite to decrease

$$\text{let } \Delta V = -\nabla C \quad \Delta C = -\|\nabla C\|^2$$

α = learn rate = how fast to descend

$$\text{let } \Delta v = -n \nabla C \quad \Delta C = -n \|\nabla C\|^2$$

the change in the cost is equal to minus eta(n) times the magnitude of the gradient ∇C

to change paras in each step

$$w_i \rightarrow w'_i = w_i - n \frac{\partial C}{\partial w_i}$$

$$b_j \rightarrow b'_j = b_j - n \frac{\partial C}{\partial b_j}$$

update rule ① to update the paras, find the gradient of the cost w/ respect to para

② times by learning rate

③ & subtract from current para value

learning rate:

- bigger = Quiver

- But too big = overshoot

- might jump to another values which isn't the minima global

Intuition \rightarrow on a mountain, identify the steepest direction & go the opposite - eventually get to the bottom

What to do when there are many paras?

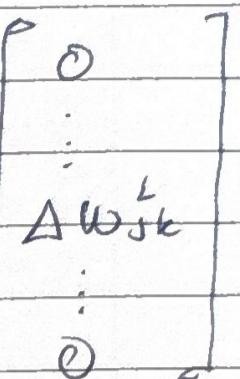
$\Delta v = -n \nabla C$ is still valid

• Scale to all layers, params & dimensions

Gradient Computation

using finite differences:

$$\frac{\partial C}{\partial w_{jk}} \approx \frac{C(w + \Delta w, b) - C(w, b)}{\Delta w_{jk}}, \quad \Delta w =$$



Change in cost over change in param

Change just one param

Problem is the number of gradient computations is equal to the number of parameters

note: $C(w, b)$ is a shared value so not repeated but the rest is

Complexity Calculation →

network w/ 784 inputs, 30 hid, 10 output

weights: $784 \times 30 + 30 \times 10 = 23,820$

Biases: $30 + 10 = 40$

each gradient comp needs 23,861 cost function calcs
per gradient

if there is many iterations of gradients

Cost per signle image $\frac{\text{Cost}}{1 \text{ gradient}}: 23,820$ (feed forward)
images = 60,000

$$23,820 \times 60,000 = 1.4 \times 10^9$$

then 1 gradient = 23,861 \star times the two for each iteration

Very big number!!!

Backpropagation

Problem w/ gradient descent is that for each iteration this is extremely expensive as it needs to keep calculating the cost function which is itself very comp expensive

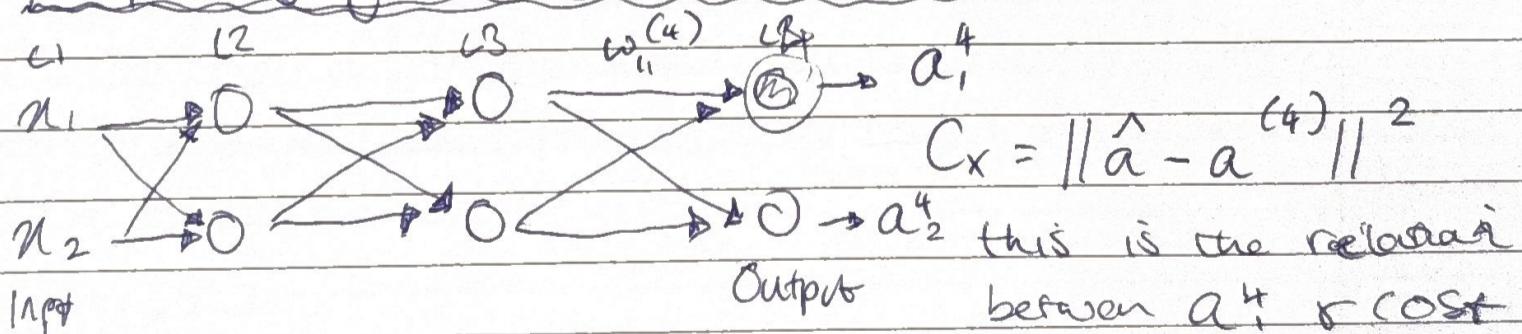
• needs to find the derivative of the cost w/ respect to each param of the cost

→ this requires computing the cost func with the new slightly change param $((w + \Delta w, b))$

Key innovation = Backpropagation

Dramatically reduces cost of gradient desc

Computing gradient using chain rule



Chain of connections:

$$C_x \leftarrow \|\hat{a} - a^4\|^2$$

$$a_1^{(4)} \leftarrow \sigma(z_1^{(4)})$$

$$z_1^{(4)} \leftarrow \#$$

$$\# = \sigma(\hat{a}_1 - a_1^{(4)})$$

$$w_{11}^{(3)} a_1^{(3)}$$

$$w_{12}^{(4)} a_2^{(3)}$$

$$w_{13}^{(4)} a_3^{(3)}$$

$$b_1^{(4)}$$

If we select a certain param we want to calculate, e.g $w_{ii}^{(4)}$, we can use the chain rule to calculate

$$\text{Chain Rule: } \frac{\partial C}{\partial w_{ii}^{(4)}} = \frac{\partial C_x}{\partial a_i^{(4)}} \cdot \frac{\partial a_i^{(4)}}{\partial z_i^{(4)}} \cdot \frac{\partial z_i^{(4)}}{\partial w_{ii}^{(4)}}$$

Here there are 3 derivatives to solve:

$$\frac{\partial C}{\partial a_i^{(4)}} = 2(\hat{a}_i - a_i^{(4)}) \quad \text{from cost}$$

$$\frac{\partial a_i^{(4)}}{\partial z_i^{(4)}} = \sigma'(z_i^{(4)}) \quad \text{from } \sigma(z_i^{(4)})$$

$$\frac{\partial z_i^{(4)}}{\partial w_{ii}^{(4)}} = a_i^{(3)} \quad \text{from } w_{ii}^{(4)} a_i^{(3)}$$

By chain rule calc =

$$\frac{\partial C_x}{\partial w_{ii}^{(4)}} = 2(\hat{a}_i - a_i^{(4)}) \sigma'(z_i^{(4)}) a_i^{(3)}$$

All of these terms are calculated using the existing activation values, instead of recalculating a cost function

$\sigma'(z_i^{(4)})$ does not just use activate but it is the derivative of the sigmoid

Deriv of sigmoid function

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

Simplified

$$\sigma'(z) = \frac{e^{-z}}{(1+e^{-z})^2} \rightarrow \sigma(z)[1 - \sigma(z)]$$

So for a sigmoid neuron:

$$a = \sigma(z)$$

a = output of neuron

$$\sigma'(z) = a(1-a)$$

So

$$\frac{\partial C_x}{\partial w_{ii}^{(4)}} = 2(\hat{a}_1 - a_1^{(4)}) \cdot a_1^{(4)}(1-a_1^{(4)}) \cdot a_1^{(3)}$$

~~Cost~~ Deriv of cost w/ respect to $w_{ii}^{(4)}$ using just values in the feedforward

Just calc'ng the product of activations that are already known

Call this $\delta_1^{(4)}$ w/out the $a_1^{(3)}$ component

$\delta_1^{(4)}$ can be reused for all $a_i^{(4)}$ terms from the neuron

$$\frac{\partial C_x}{\partial w_{ii}^{(4)}} = \delta_1^{(4)} a_1^{(3)} \quad \frac{\partial C_x}{\partial w_{i2}^{(4)}} = \delta_1^{(4)} a_2^{(3)}$$

$$\text{for the bias } \frac{\partial C_x}{\partial b^{(4)}} = \delta_1^{(4)} \cdot 1$$

* Repeat for the other output value $a_2^{(4)} \rightarrow \delta_2^{(4)}$

All computed without finite differences

Now, if we know the δ for any ~~layer~~^{layer} we can find the derivatives of the cost w/r/t to the params of that layer

$$\begin{array}{ccc} x_1 & \xrightarrow{\cancel{\delta_i^{(2)}}} & \delta_i^{(2)} \\ x_2 & \xrightarrow{\delta_2^{(2)}} & \delta_2^{(3)} \end{array} \quad \begin{array}{c} \delta_i^{(3)} \\ \delta_2^{(3)} \end{array} \quad \begin{array}{c} \delta_i^{(4)} \\ \delta_2^{(4)} \end{array} \quad \begin{array}{c} a_i^{(4)} \\ a_2^{(4)} \end{array}$$

thus for any weight or bias in net:

$$\frac{\partial C_x}{\partial w_{jk}^{(l)}} = \delta_j^{(l)} a_k^{(l-1)} \quad \frac{\partial C_x}{\partial b_j^{(l)}} = \delta_j^{(l)}$$

(8)

How to find Deltas for previous layers? $\delta^{(3)}$

Can once again use chain rule

If you know δ for output layer

& Activations of the previous layer

& the known weights of the current layer

You can find deltas of previous layer

$$\delta_j^{(4)} = \sum_k \delta_k^{(5)} w_{kj}^{(5)} a_j^{(4)} (1 - a_j^{(4)})$$

- ① In Backprop, you calc. the deltas for output layer
- ② using these deltas iteratively calc. previous deltas using just deltas & activation values
- ③ Using deltas find derivatives of all weights & bias

No finite differences, just feed forward

Gradient Descent w/ Backprop

- ① init network w/ weights & bias (FF)
- ② for each train image
 - Compute activations (FeedForward)
 - Compute δ Delta for output layer

$$\delta_j^{(4)} = 2(\hat{a}_j^{(4)} - a_j^{(4)}) a_j^{(4)} (1 - a_j^{(4)})$$
 - Compute δ Deltas for all previous layers
 - Compute gradient of cost w/ respect to the
your weight & bias for ~~each~~ the train image
- ③ Average the gradient w/ resp to each weight
& bias over the entire training set

$$\frac{\delta C}{\delta w_{jk}^L} = \frac{1}{n} \sum \frac{\delta C_x}{\delta w_{jk}^L}$$

$$\frac{\delta C}{\delta b_j^L} = \frac{1}{n} \sum \frac{\delta C_x}{\delta b_j^L}$$

- ④ update the weights & bias using grad Desc

$$w_{jk}^L \rightarrow w_{jk}^L = w_{jk}^L - \eta \frac{\delta C}{\delta w_{jk}^L}$$

$$b_j^L \rightarrow b_j^L = b_j^L - \eta \frac{\delta C}{\delta b_j^L}$$

- ⑤ repeat until cost reduces to desired level

Gradient comparsy w/ Backprop (FF)

Calc for cost of single image = 23,820 (same)

Calc for Backprop single image = 24,210

$$\text{All grads for SI: } 23,820 + 24,210 \quad (\text{sum not product}) \\ \nabla C_x = 48,030$$

$$\text{overall imgs} = 60,000 \times 48,030$$

Per iteration

$$= 10^4 \text{ improvement}$$