

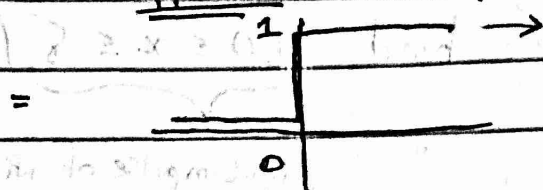
Discussion: backprop

Clarifications for HW10

Problem 1: approximation w/ ReLUs

- This is an example of universal approximators from last time (approximate any function with NN).

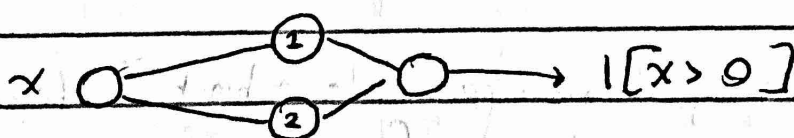
a) - You are asked to approximate a step function $1[x > 0]$



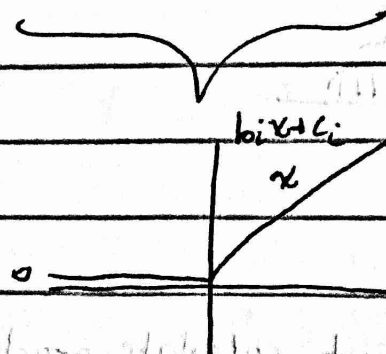
with 2 ReLU units,

$$\sum_{i=1}^2 a_i \max(0, b_i x + c_i)$$

We use ReLU as our activation function in the following network (remember from HW9):



$$= a_1 \max(0, b_1 x + c_1) + a_2 \max(0, b_2 x + c_2)$$



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

in this case $\max(0, b_1 x + c_1)$

In our case we want

$$a_1(b_1 x + c_1) + a_2(b_2 x + c_2) = 1 \quad (\text{Forget about } \phi \text{ case - ReLU takes care of that})$$

How can we choose our parameters such that this equation holds? (What is a, b, c)

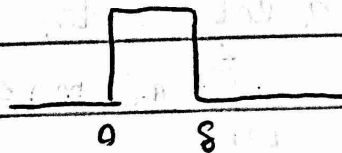
(We can see we need to subtract one ReLU from another to approximate the step function...)

* Note that it will not be perfect!

b) Using 4 ReLU's, build $1[0 \leq x \leq \delta]$

Unit impulse of width δ

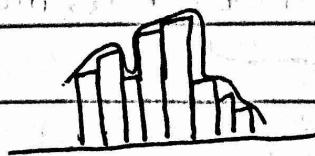
from
part a
↓



One step function can be subtracted from another (at δ) to get a box function

In this way, if we approximate a bunch of step/impulse functions, we can approximate any function

i.e.



Backpropagation

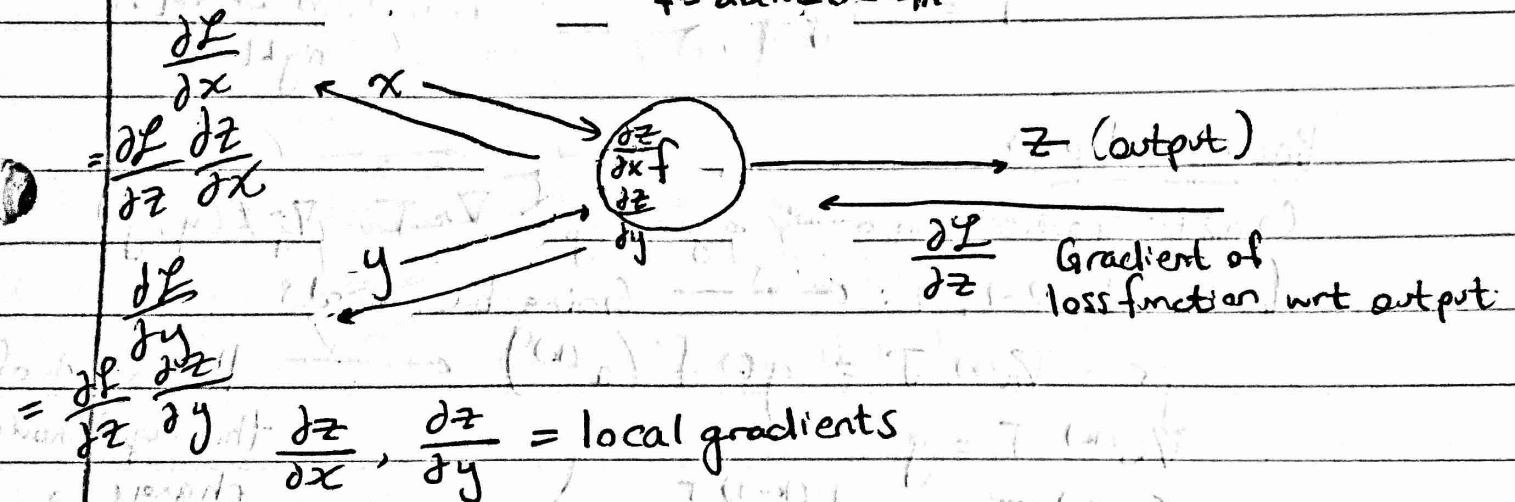
Backpropagation is the way to calculate gradients so that we can update our weights. Numerical vs. analytic gradient

$$\frac{\text{error}(w+h) - \text{error}(w)}{h} \text{ over millions of parameters} = \text{no}$$

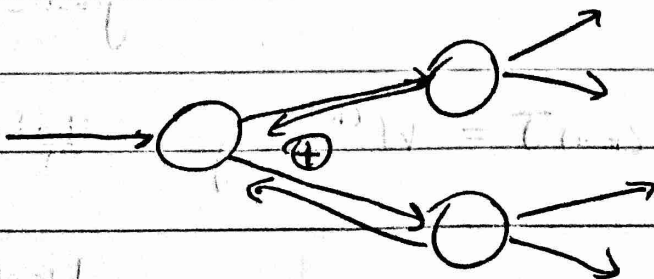
So for example, to train we would —

1. define model
 2. define error
 3. Find analytical derivative of error for every parameter w_i (backpropagation)
 4. update w_i (gradient descent)
- } any model/error that needs weight updates. not LDA, for example

Backprop example (from class) — Recursive application of the chain rule
 $f = \text{activation fn}$



Gradients will add at branches —



To do this for every node we need some general function, described in HW10.

Forward pass

$h^{(0)} = x$ "Hidden layer output" on first run

for $k = 1$ to ℓ : $\ell = \#$ of weights to update / calculate
 $a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$
 $h^{(k)} = f(a^{(k)})$
 end

$\hat{y} = h^{(\ell)}$ prediction

$J = \mathcal{L}(y, \hat{y})$ Loss (did we get it right?)

Backward pass

Compute gradient on output layer $g \leftarrow \nabla_g J = \nabla_{\hat{y}} \mathcal{L}(y, \hat{y})$

for $k = \ell, \ell-1, 1$: \leftarrow Going backwards

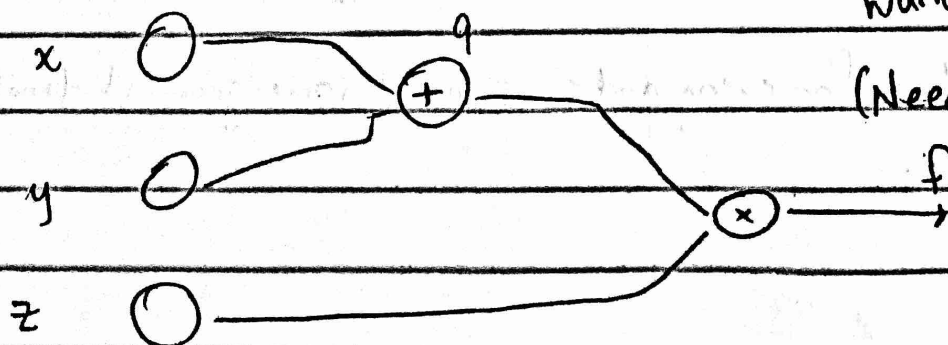
$g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$ \leftarrow How output of this layer should change to reduce error

$\nabla_{b^{(k)}} J = g$

$\nabla_{W^{(k)}} J = g \cdot h^{(k-1)T}$

We use this to get gradient of parameters

$g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$ Update



Want: $\frac{\partial \mathcal{L}}{\partial x}, \frac{\partial \mathcal{L}}{\partial y}, \frac{\partial \mathcal{L}}{\partial z}$
 (Need $\frac{\partial \mathcal{L}}{\partial g}$ to get it)