

# Deep Learning

Sunday, November 21, 2021

9:23 PM

you

$(x, f(x))$  but the nature of  $f$  is unknown.

What to try?

Halfspaces?

Kernels? Which one(s)?

Generic answer: deep Feed Forward network with nonlinear activations

Train to minimize error of output.

Gradient Descent

Loss

$$f(x) = \sum_i w_i u_i$$

$$\frac{\partial f}{\partial w_i} = u_i$$

$$f(x) = \sum_i w_i u_i \left( \sum_{i_1} w_{i, i_1} u_{i_1} \right)$$



$$\frac{\partial f}{\partial w_{i,i}} = \frac{\partial f}{\partial u_i} \cdot \frac{\partial u_i}{\partial w_{i,i}} = w_i \cdot \frac{\partial u_i}{\partial w_{i,i}}$$

$$f(x) = \sum_i w_i u_i \left( \sum_{i_1} w_{i,i_1} u_{i_1} \left( \sum_{i_2} w_{i_1,i_2} u_{i_2} \right) \right)$$



$$\frac{\partial f}{\partial w_{i_2}} = \sum_i \frac{\partial f}{\partial u_i} \cdot \frac{\partial u_i}{\partial w_{i_2}}$$

$$= \sum_i w_i \cdot \frac{\partial u_i}{\partial u_{i_1}} \cdot \frac{\partial u_{i_1}}{\partial w_{i_2}}$$

$$= \left( \sum_i w_i \cdot w_{i,i_1} \right) \frac{\partial u_{i_1}}{\partial w_{i_2}}$$

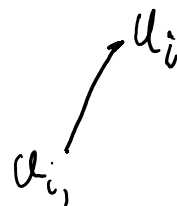
Computing  $\frac{\partial f}{\partial u}$  for all  $u$  suffices

Backprop.

-  $\frac{\partial f}{\partial f} = 1$ , send downward

For node  $u$

- sum all incoming values to get  $B$
- send  $B \cdot \frac{\partial u}{\partial u_{i_1}}$  to each  $u_{i_1}$  below.



Lemma.  $\delta = \frac{\partial f}{\partial u}$

Pf. Induction + chain rule.

base:  $\frac{\partial f}{\partial f} = 1$ .

step:  $u_{i_1}$  receives  $\frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial u_{i_1}}$  from  $u$

Hence computes induction hypothesis

$$\sum_u \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial u_{i_1}} = \frac{\partial f}{\partial u_{i_1}}$$

Time:  $O(M)$ .

Fast, general, but is it any good?

— If no hidden layer then  $GD \rightarrow OPT$ .

— If hidden layers, heavily overparametrized

more parameters than data  
then  $GD \rightarrow OPT$ . (But does it generalize?!)

Thm. ... approximate any

Thm 1. Neural Networks can approximate any continuous function. [Cybenko 89, Hornik-Stinchcombe-White 89, Barron 93]

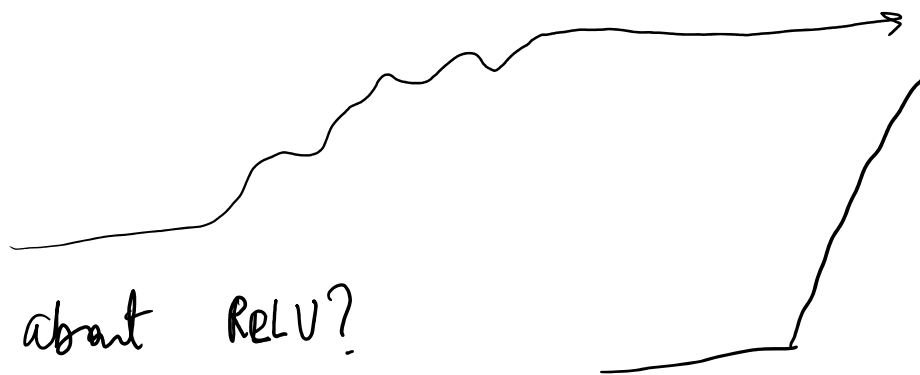
In fact, depth 2 suffices! Sufficiently wide (= 1 hidden layer)

But depth can help —  $\exists$  functions that need a large representation with small depth and a small representation with moderate depth.

Thm 1 needs activation units to be "sigmoidal"

$$\sigma(x) \rightarrow 0 \quad x \rightarrow -\infty$$

$$\sigma(x) \rightarrow 1 \quad x \rightarrow \infty$$



What about ReLU?

Put two together:

\*

put two together:

