

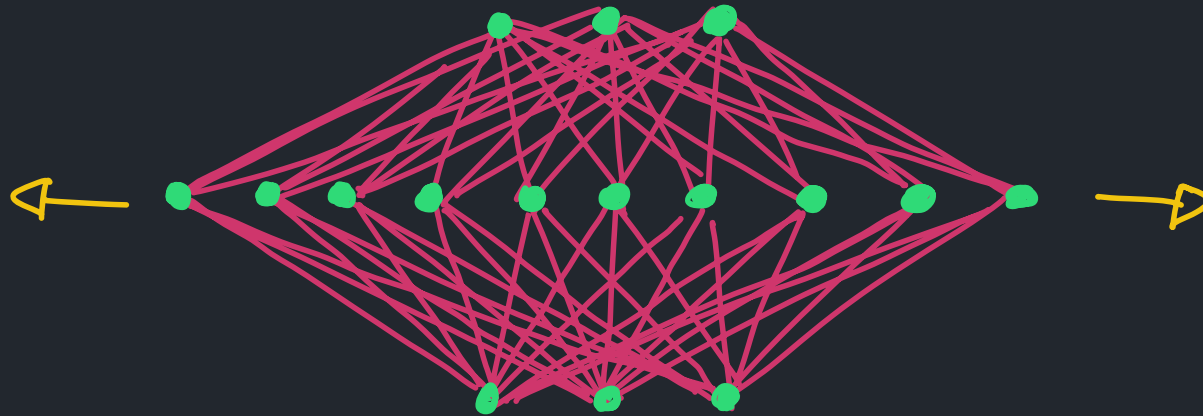
# Approximation Theory

---

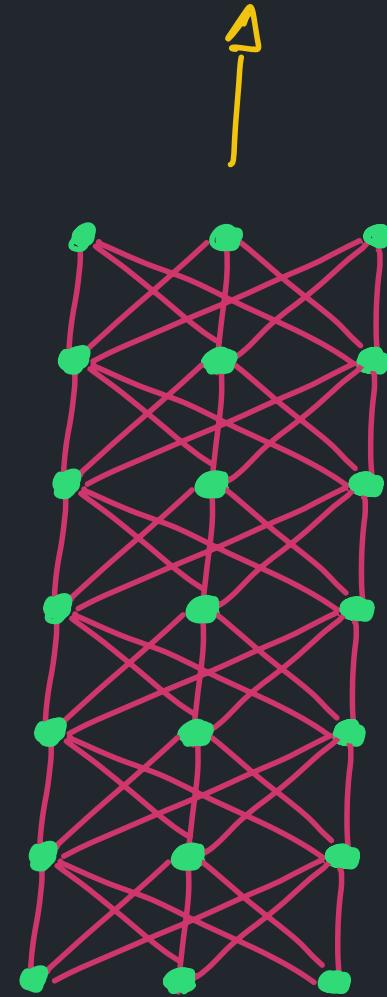
Jeremy Bernstein  
jbernstein@mit.edub



Would you rather...



scale width



scale depth

# Today's Lecture

What class of functions can a neural net express?

"neural nets are universal function approximators"

Then what architecture should I use?

"three layers are enough"

"stack more layers!!!"

...then what should I do in practice?

# The machine learning puzzle

Three pieces to the puzzle:

- ① Approximation Does there exist a neural net in my model family that fits the training data?
- ② Optimization If it does exist, can I find it?
- ③ Generalization Does it work well on unseen data?

This lecture will focus mainly on the first question.

# A motivating problem



Can we classify this data with the function:

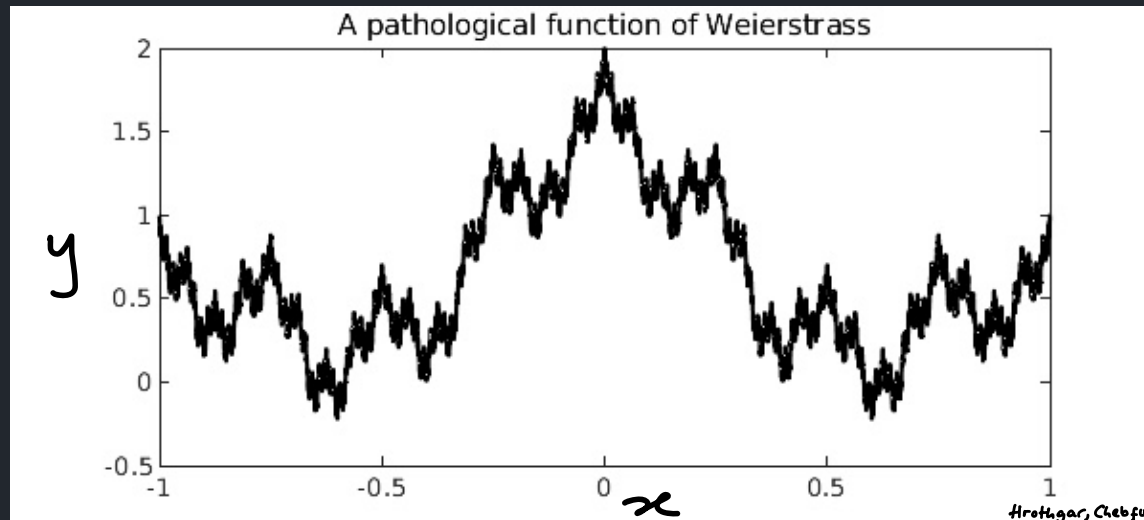
$$f(x) = \text{relu}(w^T x + b) \quad ?$$

Or what about with "two layers":

$$f(x) = \sum_i \alpha_i \text{relu}(w_i^T x + b_i) + \beta_i \quad ?$$

# Another motivating problem

An example of a "fractal curve"



Weierstrass' function is everywhere continuous  
nowhere differentiable

→ Can you fit it with a neural network?

If so, how big should the network be?

# Formalizing the approximation problem

Given a family of curves  $G$  — exclude pathological functions

And a family of neural networks  $F$  — e.g. 5 layer relu MLPs

For any curve  $g \in G$

Does there exist a neural net  $f \in F$  think: small number

Such that  $\text{error}(f, g) < \varepsilon$  ?

e.g.  $\text{error}(f, g) \triangleq \max_x |f(x) - g(x)|$  " $L_\infty$  error"

or  $\text{error}(f, g) \triangleq \int dx |f(x) - g(x)|$  " $L_1$  error"

# One nice family of curves $G$

## Lipschitz continuous functions

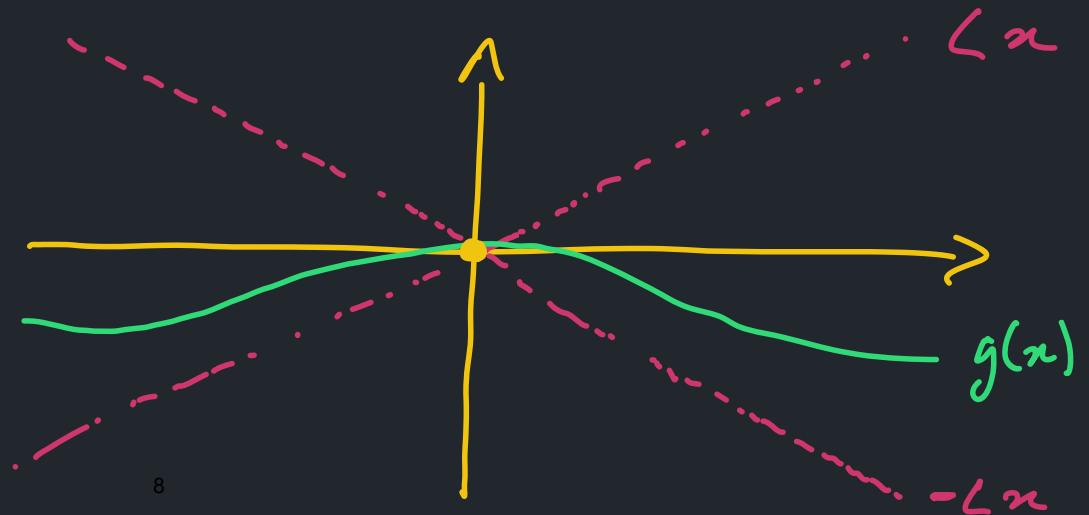
$g: \mathbb{R} \rightarrow \mathbb{R}$  is " $L$ -Lipschitz" if

$$|g(x + \Delta x) - g(x)| \leq L |\Delta x|$$

for all inputs  $x \in \mathbb{R}$  and all  $\Delta x \in \mathbb{R}$ .

Intuition: the slope of  $g$  cannot exceed  $L$ .

If  $g$  passes through the origin then it can never stray outside  $\pm Lx$ .





# Extending to multi-dimensional inputs

## Lipschitz continuous functions

$g: \mathbb{R}^d \rightarrow \mathbb{R}$  is " $L$ -Lipschitz" if

$$|g(x + \Delta x) - g(x)| \leq L \|\Delta x\|_{\text{RMS}}$$

for all inputs  $x \in \mathbb{R}^d$  and all  $\Delta x \in \mathbb{R}^d$ .

Define the "RMS-norm",

$$\|x\|_{\text{RMS}} \triangleq \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}$$

$L$ , it measures the average (root-mean-square) size of the entries of the vector.

In this lecture, we will prove...

d-dimensional hypercube.

Theorem:

Let  $g: [0,1]^d \rightarrow \mathbb{R}$  be any  $L$ -Lipschitz function.

Then for any error  $\varepsilon > 0$

There exists a 3-layer relu network

With  $N = 4d(L/\varepsilon)^d$  units

such that  $\int_{[0,1]^d} |f(x) - g(x)| dx < 2\varepsilon$

# Strategy for proving the result

## Step one:

derive result for 1d inputs  
via approximation with  
"rectangular strips"

ie ignore relu  
networks for now.

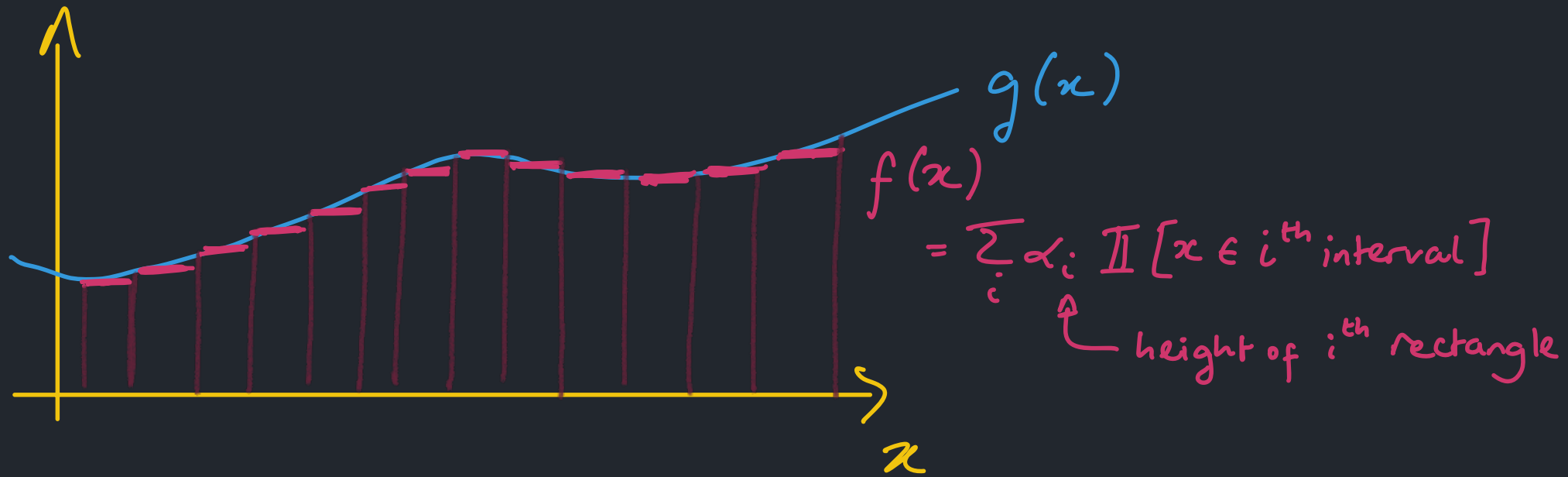
## Step two:

generalise to higher input  
dimension

## Step three:

Show that relu networks can  
approximate rectangular strips

# Step One: Approximation with rectangles



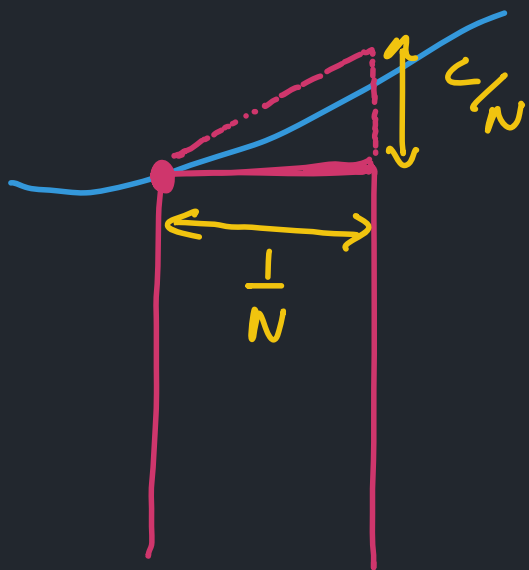
Given  $N$  rectangular strips, what is the approximation error?

CLAIM:  $\text{error} \leq \frac{L}{2N}$

increases with  
Lipschitz constant

decreases with  
number of strips

$N$  strips  $\Rightarrow$  each strip has width  $\frac{1}{N}$



by Lipschitzness the curve can't exceed the top side of the triangle

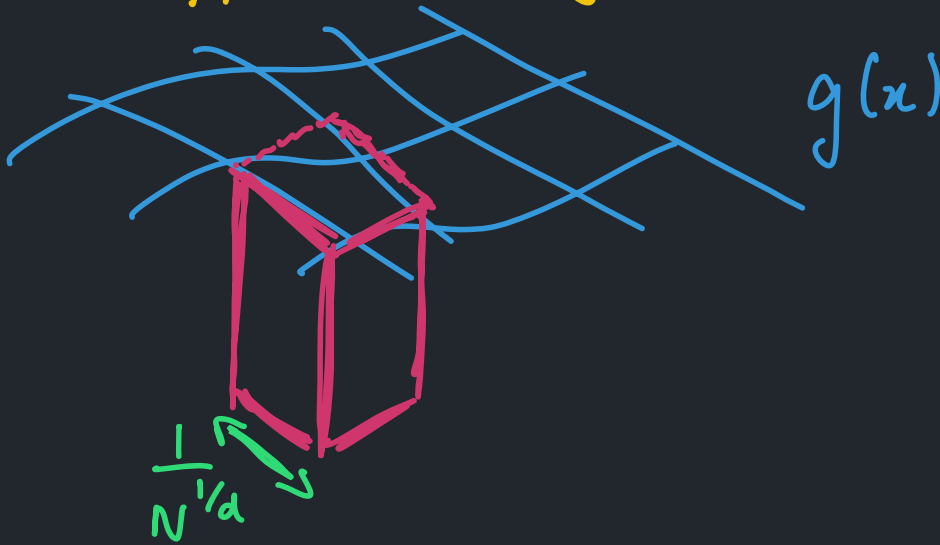
The triangles has area  $\frac{1}{2} \frac{L}{N^2}$

$$\begin{aligned} \text{Total error} &= \int dx |f(x) - g(x)| \leq N \times \frac{1}{2} \frac{L}{N^2} \\ &= \frac{1}{2} \frac{L}{N} \end{aligned}$$

So to achieve error  $\varepsilon$  we need  $N = \frac{1}{2} \frac{L}{\varepsilon}$

Step two: Higher input dimension

Think: approximating a surface with cuboids



$N$  hyperrectangles  $\Rightarrow$  each has side of Euclidean length  $\frac{1}{N^{1/d}}$

$$\text{total error} = \int dx |f(x) - g(x)| \leq N \times \frac{L}{N^{1/d}} \times \frac{1}{N} = \frac{L}{N^{1/d}}$$

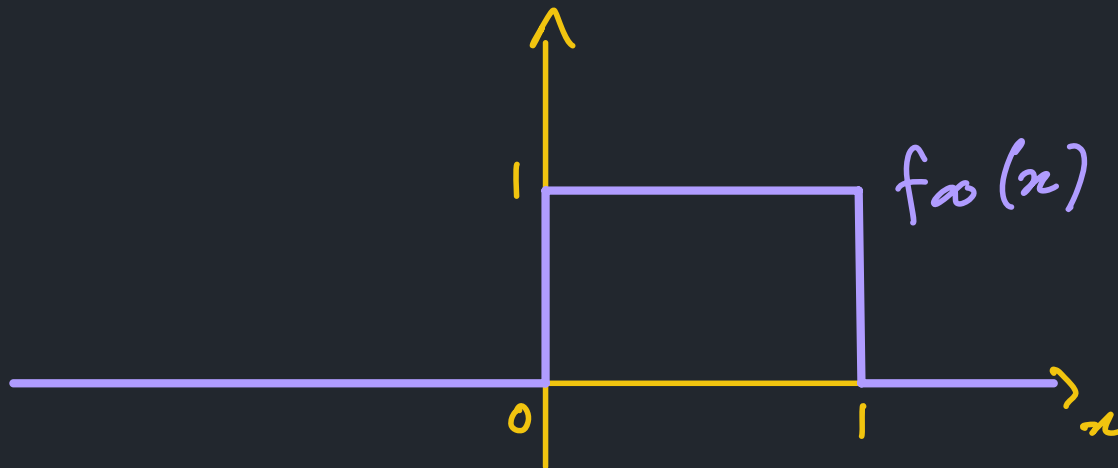
number of hyperrectangles      height of error cap      area of error cap.

To get error  $\epsilon$  requires  $N = \left(\frac{L}{\epsilon}\right)^d$  hyperrectangles

# Step Three: Relu networks can fit rectangles

Now define  $f_c(x) = \begin{bmatrix} +1 \\ -1 \\ -1 \\ +1 \end{bmatrix}^T \text{relu} \begin{bmatrix} cx \\ cx - 1 \\ c(x-1) - 2 \\ c(x-1) - 3 \end{bmatrix}$

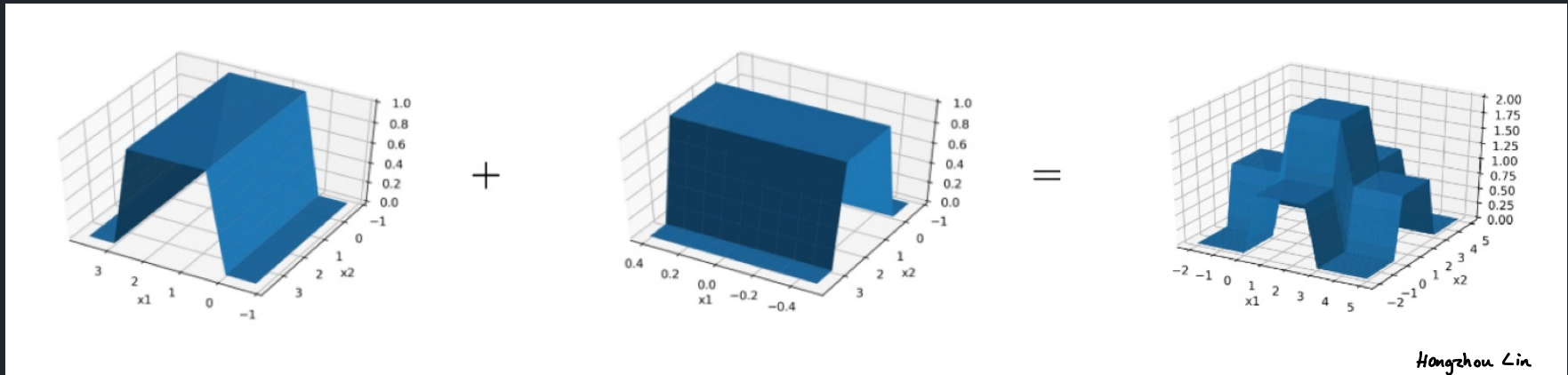
CLAIM: let  $c \rightarrow \infty$  and we get:



↳ can also translate horizontally by adjusting weights and biases.

But we need  $d$ -dimensional hyperrectangles...

Solve by adding 1-dimensional rectangles and thresholding appropriately



Only exceeds  $d-1$  when all rectangles are "on"


↳ so just threshold at  $d-1$ .



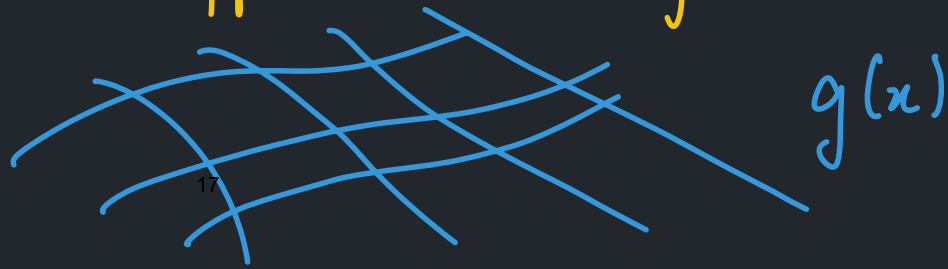
# Assembling the pieces

rectangle:  $f_c(x) = \begin{bmatrix} +1 \\ -1 \\ -1 \\ +1 \end{bmatrix}^T \text{relu} \begin{bmatrix} cx \\ cx - 1 \\ c(x-1) - 2 \\ c(x-1) - 3 \end{bmatrix}$

hyperrectangle:  $h_c(x) = \text{relu} \left[ \sum_{i=1}^d f_c(x_i) - (d-1) \right]$

linear combination of hyperrectangles  $f(x) = \sum_i \alpha_i h_c(x - u_i)$    
  position of i-th gridpoint

Let  $c \rightarrow \infty$  and we get an approximation of an arbitrary Lipschitz surface!



# Some comments

Theorem:

Let  $g: [0,1]^d \rightarrow \mathbb{R}$  be any  $L$ -Lipschitz function.

Then for any error  $\varepsilon > 0$

There exists a 3-layer relu network

With  $N = 4d(L/\varepsilon)^d$  units

such that  $\int_{[0,1]^d} |f(x) - g(x)| dx < 2\varepsilon$

General idea: approximate "bumps" then linearly combine

Needs exponentially many neurons in dimension

Taking  $c \rightarrow \infty$  feels unrealistic

Approximating rectangles feels like a trick

# Imagine training the rectangle representation



- would just get rectangles on the training points  $x$
- regularisation (weight decay) would suppress the others
- would not generalise!

# Further reading

More results on universal function approximation

- Barron's theorem

"Smooth functions can be approximated with fewer neurons"  
leverages Fourier representation

- 2 layers are enough

e.g. Hornik, Stinchcombe and White (1989)

uses Stone-Weierstrass theorem

Is universal function approximation important?

Is "UFA" sufficient for learning to work?

No, there are many UFAs that we usually don't do ML with:

Examples:      Fourier series

                 polynomials

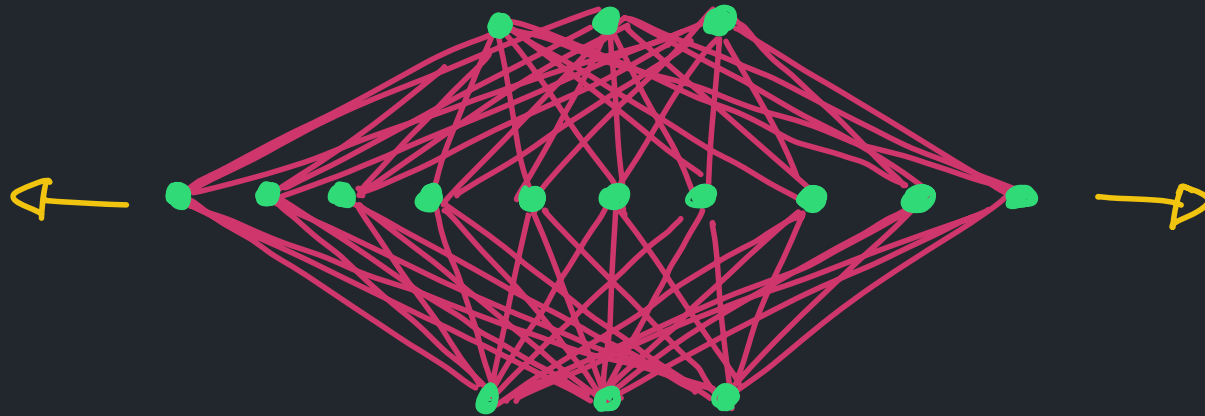
                 the space of Python programs

Is "UFA" necessary for learning to work?

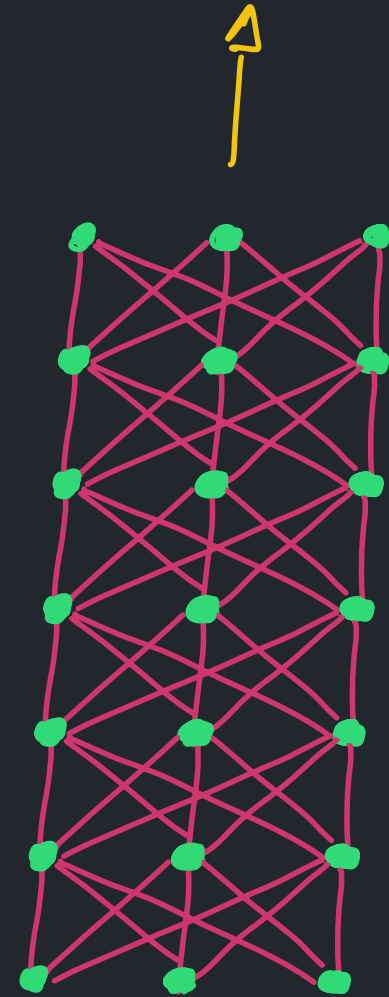
Width versus depth

---

Would you rather...



scale width



scale depth

# Width versus depth

## Advantages of scaling width

- 3 layer (or even 2 layer) NNs are "universal function approximators"
- width is inherently parallelisable, depth is sequential
- width is easier to train, depth leads to "compound problems"

So, scaling width is obviously better!



## Or is it?.... Depth separations

Universal function approximation results suggest needing exponentially many hidden units at small width

"depth separation" results construct deep networks that require exponentially more units to fit with a shallow network.

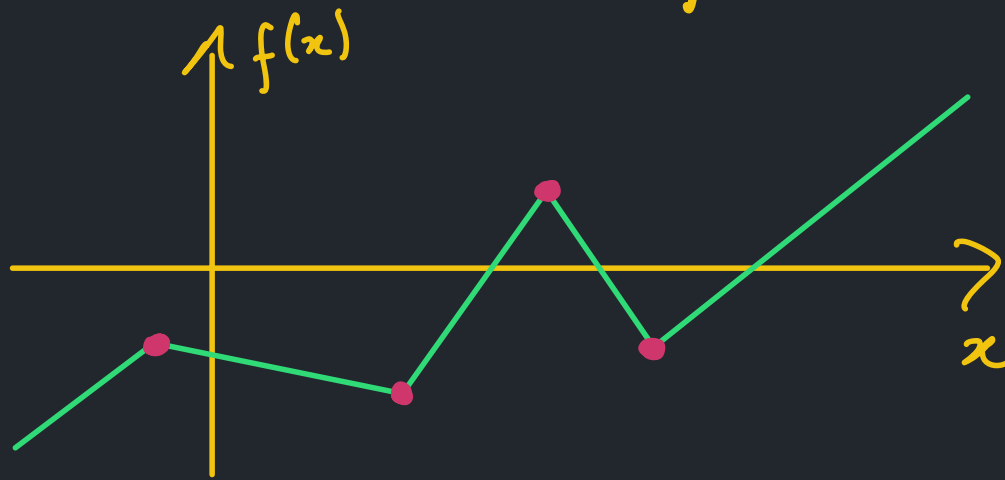
# The shape of a result

To prove a "depth separation"

e.g. "number of linear regions"

- ① pick a property of a function
- ② construct a deep network that has this property
- ③ prove that a shallow network would need exponentially more units to also have this property

# Piecewise linear functions



we will define a "kink"  
to be a place where the  
gradient changes

Here, #kinks = 4

Claim: relu networks are piecewise linear ("PWL")

Why? relu is PWL

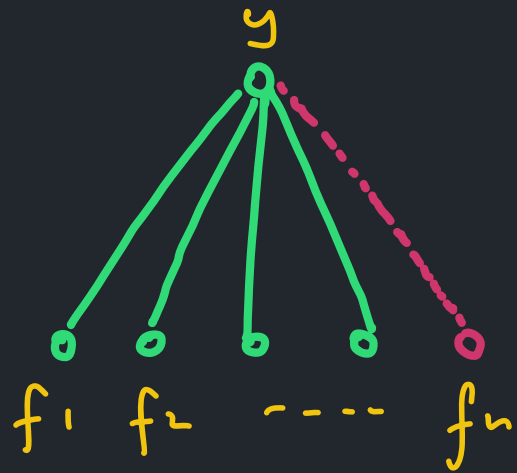
$f, g$  both PWL  $\Rightarrow f + g$  PWL

$f, g$  both PWL  $\Rightarrow f \circ g$  PWL

$f$  PWL  $\Rightarrow \alpha \cdot f$  PWL for  $\alpha \in \mathbb{R}$

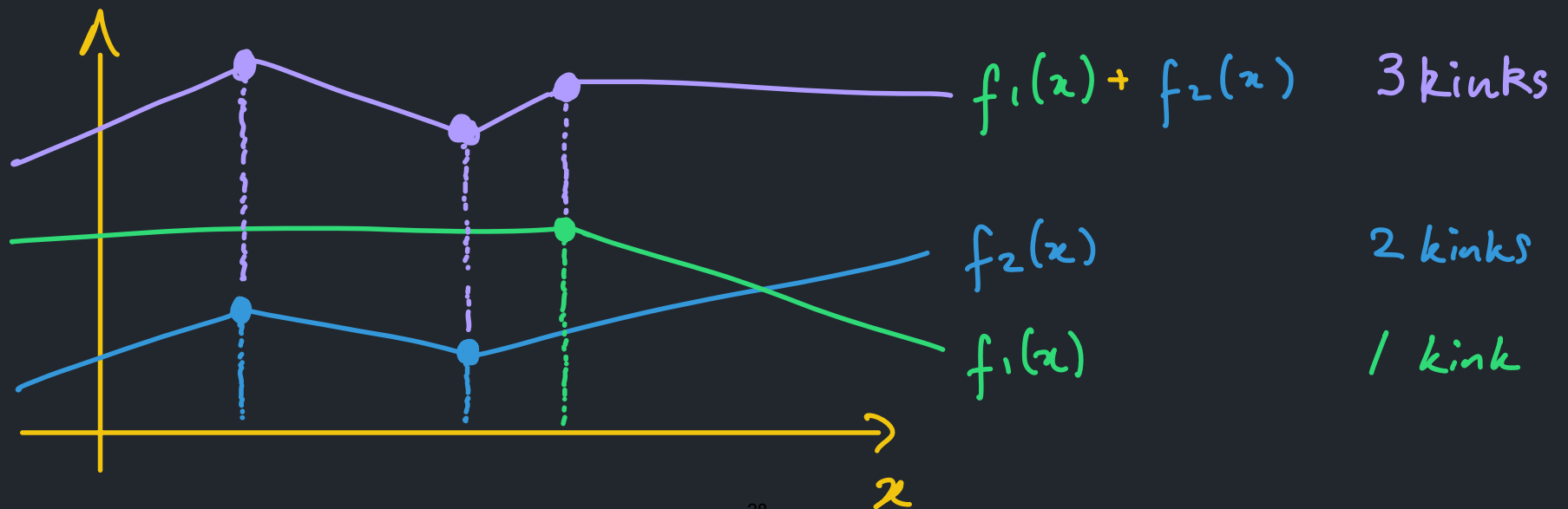
Okay, let's construct a depth<sup>27</sup> separation based on #kinks

# Intuition: Effect of adding width

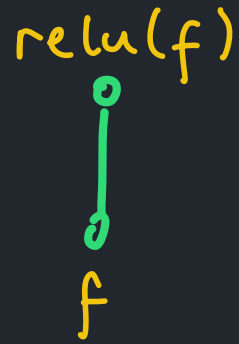


$$y(x) = \sum_{i=1}^n \alpha_i f_i(x)$$

When we add functions, at most we add the # kinks.

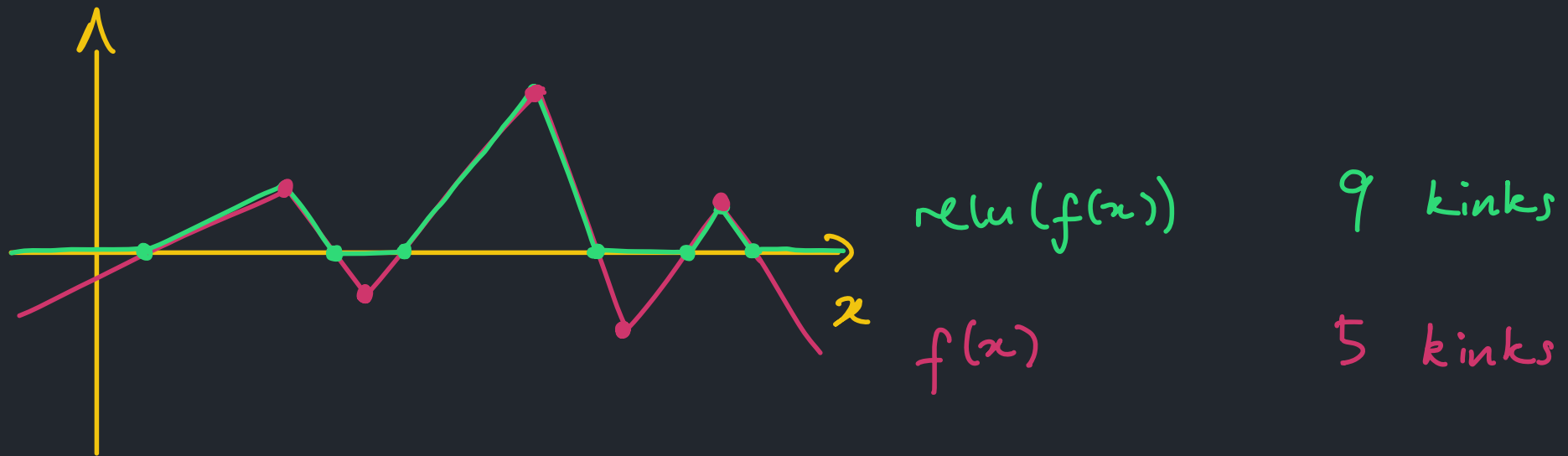


# Intuition: Effect of applying relu



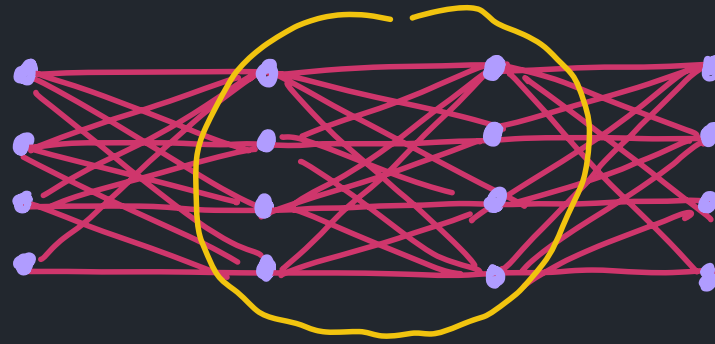
$$y(x) = \text{relu}(f(x))$$

When we apply relu, at most we double the # kinks



Because linear pieces can be split in two.

More formally



Consider the  $l^{\text{th}}$  layer:

$$f_l(x) = \text{relu} \left( W_l f_{l-1}(x) + b_l \right)$$

Diagram illustrating the equation above:

- $x$  is labeled "vector in  $\mathbb{R}^n$ " with an arrow pointing to it.
- $W_l$  is labeled " $n \times n$  matrix" with an arrow pointing to it.
- $f_{l-1}(x)$  is labeled "vectors in  $\mathbb{R}^n$ " with an arrow pointing to it.
- $b_l$  is labeled "vectors in  $\mathbb{R}^n$ " with an arrow pointing to it.

Let  $\text{KINKS}_l$  denote the max number of kinks over the  $n$  coordinates of  $f_l(x)$ .

Then it holds that  $\text{KINKS}_l \leq 2n \cdot \text{KINKS}_{l-1}$

Since  $\text{KINKS}_0 = 1$ , this implies  $\boxed{\text{KINKS}_L \leq (2n)^L}$

# Interpreting the result

We showed that:  $\boxed{KINKS_L \leq (2n)^L}$

$KINKS_L$  = # kinks in the function at layer  $L$

$n$  = width

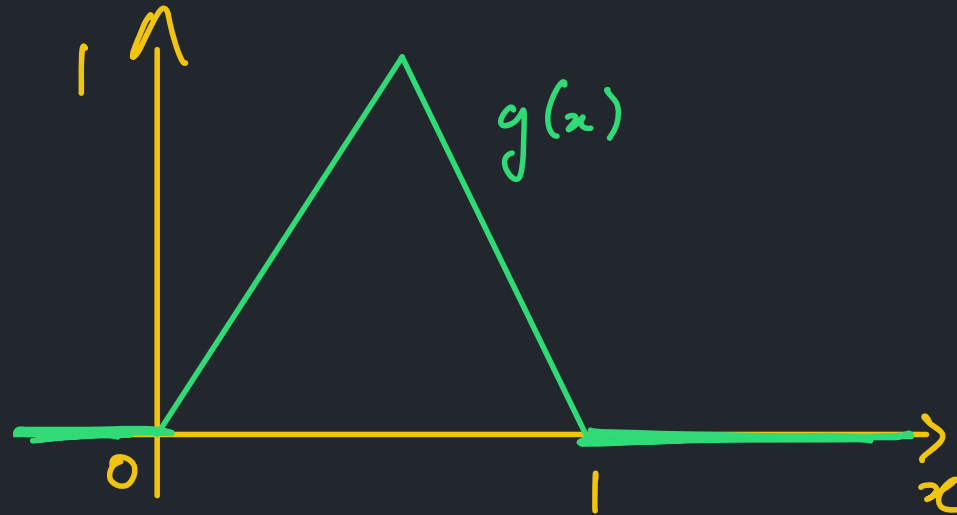
$L$  = depth

↳ upper bound grows at best polynomially in width  
but exponentially in depth

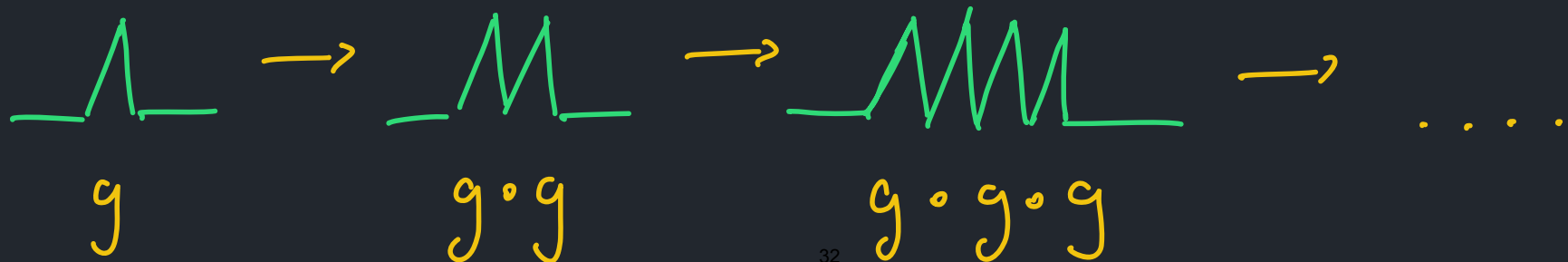
Need to ask: is the bound ever attained?

Short answer: Yes!

Define  $g(x) = \text{relu} \left[ 2 \cdot \text{relu}(x) - 4 \cdot \text{relu}(x - \frac{1}{2}) \right]$



Each time we "iterate"  $g$  — i.e. compose with self — we double the number of linear regions





## $\Rightarrow$ Depth separation

- $g \circ g \circ \dots \circ g$  500 times has  $2^{500} - 1$  kinks
- it has 1000 layers of width 2
- to get the same number of kinks with a 3-layer network we would need a width  $n$  of

$$\text{KINKS} \leq (2n)^L$$

$$\Rightarrow n \geq \frac{1}{2} \cdot \text{KINKS}^{\frac{1}{L}}$$

$$= \frac{1}{2} \cdot (2^{500} - 1)^{\frac{1}{3}}$$

$$= \boxed{7 \times 10^{49} \text{ units}}$$

# What does this not say?

It does not mean that :

- very deep networks are easy to train  
OPTIMISATION

- very deep networks would generalise well  
GENERALISATION

In our machine learning puzzle, it only tells us something about APPROXIMATION!

## Further reading

- Telgarsky (2015, 2016)
  - our depth separation
- Safran and Shamir (2017)
  - a different depth separation
- Lu, Pu, Wang, Hu and Wang (2017)
  - results on "minimum width" needed to be a universal function approximator even at "large depth"
  - relates to rank of the weight matrices

## Practical considerations

---

Let's consider the whole puzzle

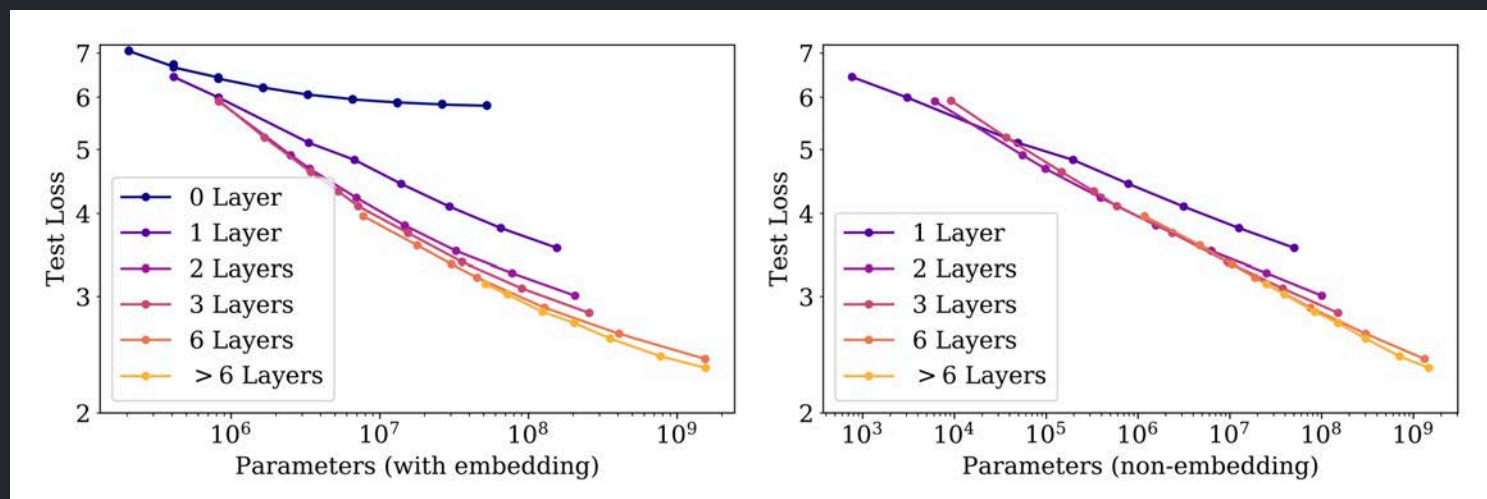
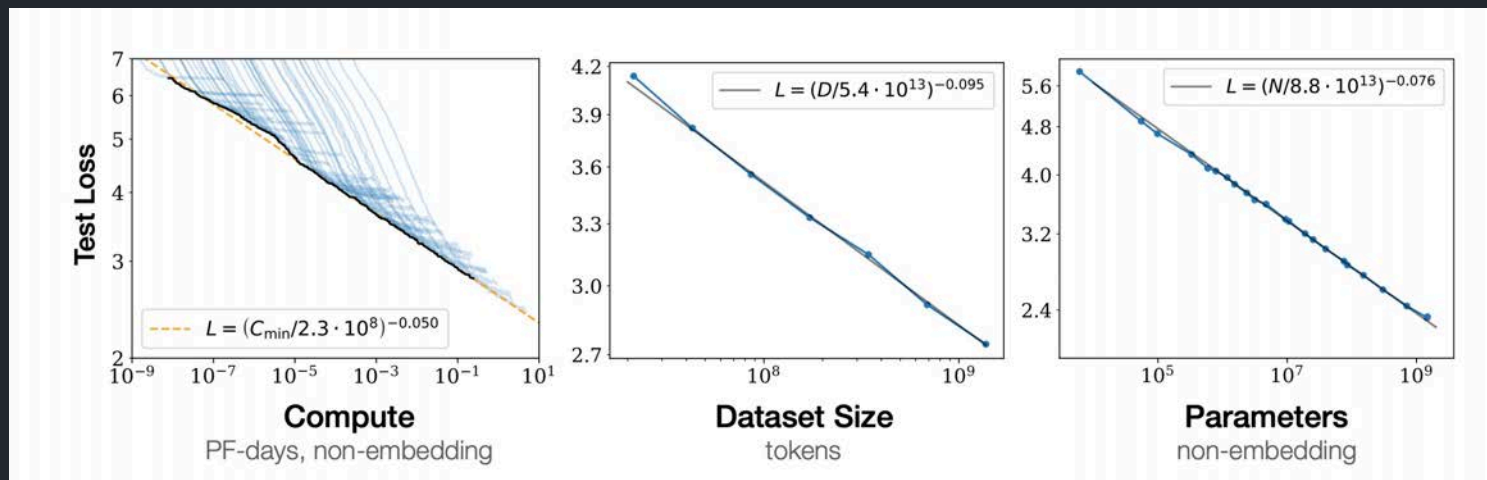
APPROXIMATION

+ OPTIMISATION

+ GENERALISATION

- Pretend you work at an LLM startup
- You want to train the most efficient LLM possible
- You really care about the optimal width vs. depth

# Scaling laws



From Kaplan, McCandlish et al (2020)

# The importance of confounders

## "Chinchilla scaling rules"

Hoffmann, Borgeau, Mensch et al (2020)

- question some results in Kaplan et al
- suggest a different learning rate schedule
- changes certain scaling results

—> to really answer

"what is the optimal width versus depth"  
need to obsess over "minor details" of the  
training pipeline

Wrapping up

---



# Summary

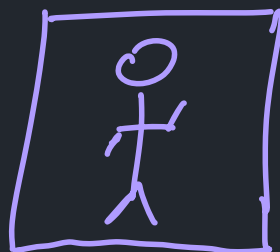
- Very wide shallow neural nets are universal function approximators
- Deeper networks can fit certain kinds of function with many fewer neurons — "depth separations"
- Unclear how these results interact with training and generalisation

# Preview: Inductive biases

Suppose we want to solve two different machine learning problems



"hello"



"human"

In both cases, we could flatten the inputs into vectors and just apply an MLP...

↳ Hey, it's a universal function approximator.  
But, is this a good idea?

MIT OpenCourseWare

<https://ocw.mit.edu>

6.7960 Deep Learning

Fall 2024

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>