# Network Function Spaces



Jeremy Bernstein
[bernstein@caltech.edu](mailto:bernstein@caltech.edu)
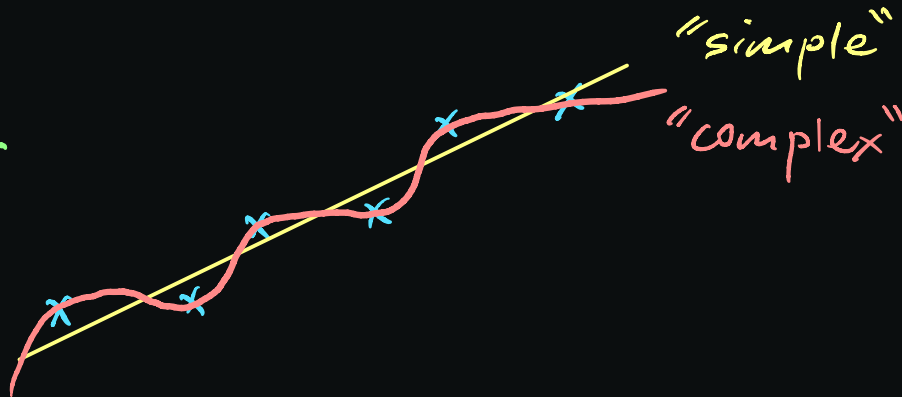
# Agenda for today

1. Complexity of a function space
2. Universal function approximators?
3. Studying random functions
4. Neural networks as Gaussian processes

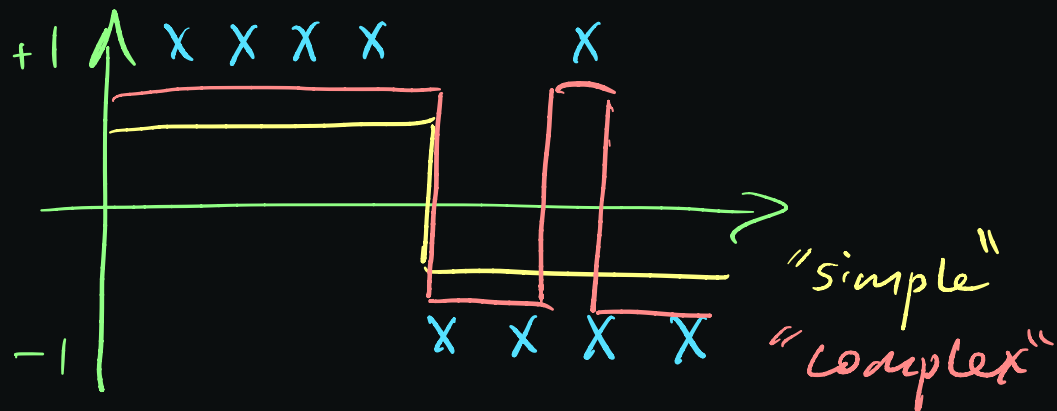# A basic question

- What does an NN's function space look like?

How "complex" are the functions?

Regression

"simple"
"complex"

Classification

+1

x x x x          x

"simple"

−1              "complex"

here we are being vague about what "complex" means

3

# Universal function approximation

- Theorems like:

    *"a wide enough NN can fit any function"*

- Empirical findings like:

    *"my NN can fit any labelling of the train set"*

But this does not address what kinds of function
are common —— what kind of function is the
architecture biased toward?

a.k.a. the NN's inductive bias —— the focus of this lecture.

# Weight space $\mapsto$ function space

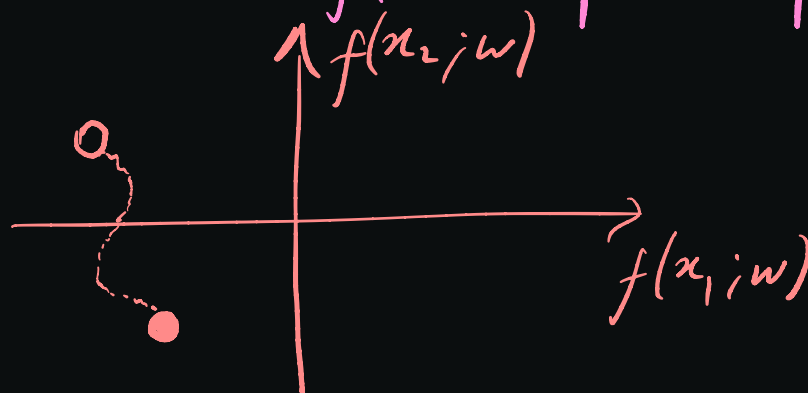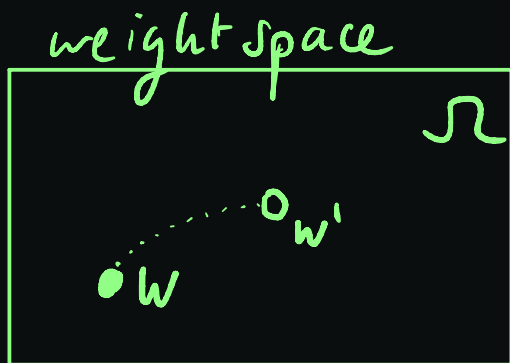Assuming $n$ inputs and an NN with a 1D output,

weight vector
$w$

$\mapsto$

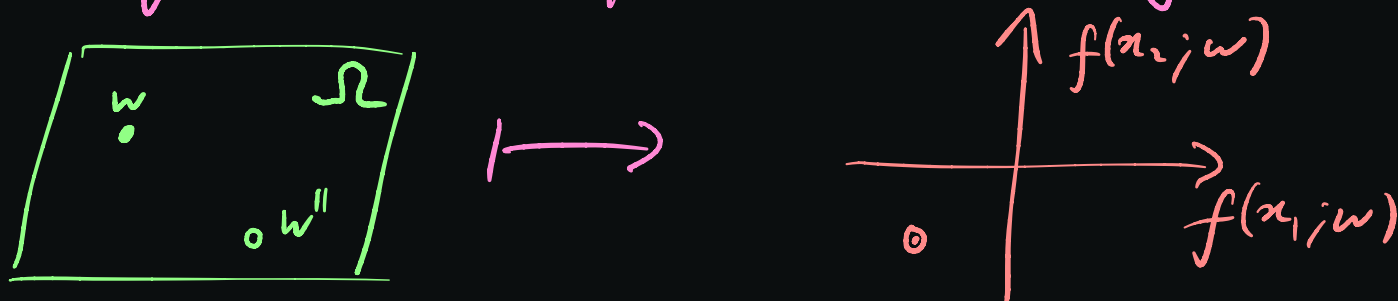function space representation

$$\begin{bmatrix} f(x_1; w) \\ f(x_2; w) \\ \vdots \\ f(x_n; w) \end{bmatrix} \in \mathbb{R}^n$$

Varying $w$ in weight space also moves the function space rep.
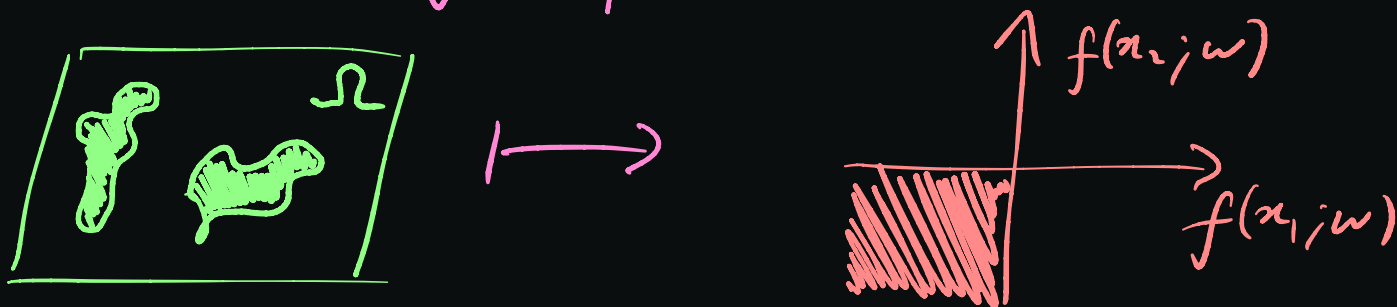
weight space



$\mapsto$

# Life is simpler in function space

① Many weight vectors map to the same function:



② Simple geometries in function space may be intractably complicated in weight space:
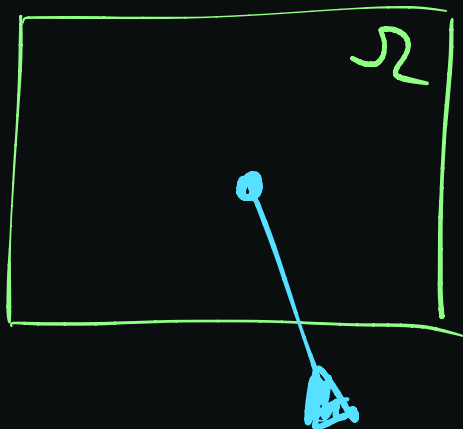


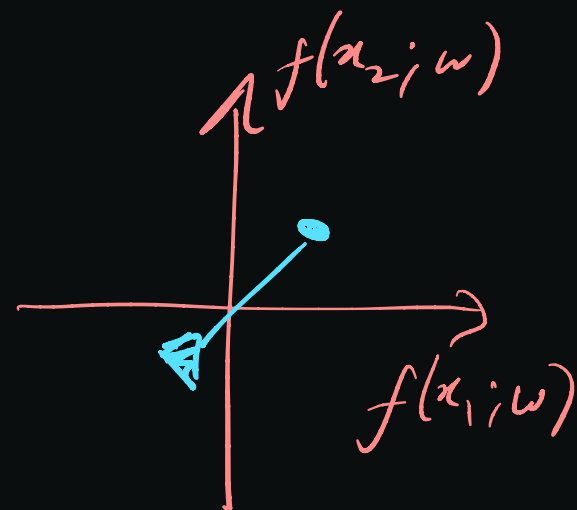...so how do we move there?

# Random functions

Key idea: to study what kinds of function an NN is biased towards, study the properties of random functions that the NN implements.

# Sampling functions

Imagine throwing a dart at weight space...

and inspecting the corresponding fn.

... then throw more darts

$f(x_2; w)$

$f(x_1; w)$

# An example

network details $\left\{\begin{array}{l}\text{3 layer MLP} \\ \text{width 1000} \\ \text{1 output}\end{array}\right.$

Repeatedly sample the weights randomly.

Input 2     Input 3     Input 4



Input 1



Outputs jointly Gaussian!

Jupyter notebook on the course website.

# Output correlations

A Gaussian with mean zero is fully specified by its covariance matrix.

When sampling random networks, this is given by

$$\Sigma_{ij} := \mathbb{E}_{w \sim \mathbb{P}}\left[ f(x_i; w)\, f(x_j; w) \right]$$

measure on weight space

output on input $i$

output on input $j$

# Wide NNs are Gaussian processes

*enough* (written above "Wide")

- For any finite collection of inputs $x_1, x_2, \ldots, x_n$.

- For weights $w \overset{iid}{\sim}$ Gaussian.

- For a wide enough NN.

— the outputs are jointly Gaussian.

$$\begin{bmatrix} f(x_1; w) \\ \vdots \\ f(x_n; w) \end{bmatrix} \sim \mathcal{N}(\mu, \Sigma).$$

# Aside: Gaussian processes

If for all finite collections of inputs

$$x^{(1)}, ..., x^{(n)}$$

the following holds:

$$f(x^{(1)}), ..., f(x^{(n)}) \sim \mathcal{N}(\mu, \Sigma)$$

then we say that $f$ is a <u>Gaussian process</u>.

To prove the claim that random, sufficiently wide NNs behave like GPs, we will need some...

# Gaussian facts !

# **Classical CLT**

*central limit theorem*

Let $X_1, \ldots, X_n$ be i.i.d. random variables each with mean 0 and finite variance $\sigma^2$.

Then as $n \to \infty$,

$$\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \longrightarrow 0 \qquad \text{with probability one.}$$

$$\sqrt{n} \cdot \overline{X} = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} X_i \longrightarrow \mathcal{N}(0, \sigma^2).$$

14

# Multivariate CLT

Let $Y_1, \ldots, Y_n$ be i.i.d. random vectors each with mean 0 and finite covariance $\Sigma$.

Then as $n \to \infty$,

$$\overline{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i \longrightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{with probability one.}$$

$$\sqrt{n} \cdot \overline{Y} = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} Y_i \longrightarrow \mathcal{N}(0, \Sigma).$$

# Linear transformation of Gaussian

Let $Y$ be a random vector $Y \sim \mathcal{N}(\mu, \Sigma)$.

What is the distribution of $AY$?

$\hookleftarrow$ fixed matrix.

$AY$ is also Gaussian, by moment generating fns.

To see this,
- MGF of $Y$ is: $\mathbb{E} e^{t^T Y} = e^{t^T \mu + \frac{1}{2} t^T \Sigma t}$
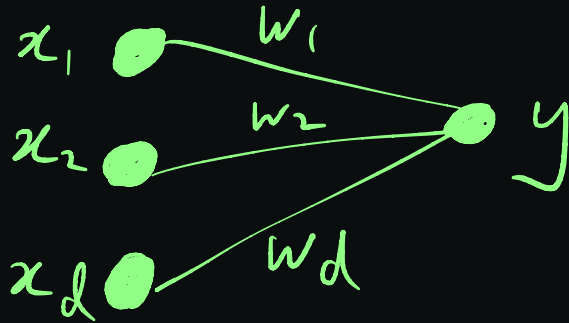
$\Rightarrow$ MGF of $AY$ is: $\mathbb{E} e^{t^T AY} = e^{t^T (A\mu) + \frac{1}{2} t^T A \Sigma A^T t}$

$\Rightarrow AY \sim \mathcal{N}(A\mu, A\Sigma A^T)$ $\left( \begin{array}{c} \text{by uniqueness of} \\ \text{MGFs.} \end{array} \right)$

16

# NNGP correspondence

- We now have everything we need to prove that NNs behave like GPs.

- We'll start simple and build up.

# Linear neuron

$x_1$ ●━━━ $w_1$

$x_2$ ● $w_2$ ● $y$

$x_d$ ● $w_d$

Consider $n$ inputs:

$$x^{(1)}, \ldots, x^{(n)} \in \mathbb{R}^d$$

Stack the inputs into a "data matrix":

$$X = \begin{bmatrix} \underline{\quad} x^{(1)} \underline{\quad} \\ \underline{\quad} x^{(2)} \underline{\quad} \\ \vdots \\ \underline{\quad} x^{(n)} \underline{\quad} \end{bmatrix}$$

Then the $n$ outputs may be written as $Xw$

If the components of $w$ are jointly Gaussian, then the outputs are given by a linear transformation of $w$ ⟹ the outputs are jointly Gaussian.

18

# Linear neuron: covariance

What is the covariance of the linear neuron GP?

Consider two outputs

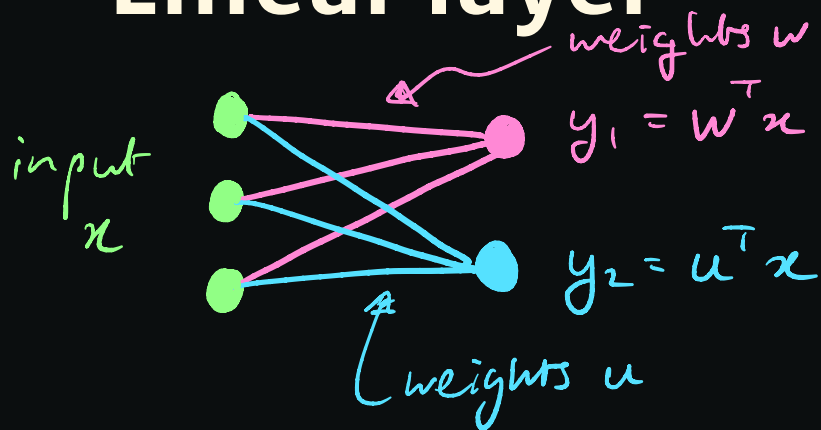$$y^{(1)} = \sum_{i=1}^{d} w_i x_i^{(1)} \qquad\qquad y^{(2)} = \sum_{i=1}^{d} w_i x_i^{(2)}$$

Suppose the weights $w_i \overset{iid}{\sim} \mathcal{N}(0, \sigma^2)$.

The means are $\mathbb{E}\, y^{(1)} = \mathbb{E}\, y^{(2)} = 0.$

Then the covariance is

$$\mathbb{E}\left[ y^{(1)} y^{(2)} \right] = \mathbb{E}\left[ \sum_{i,j} w_i^{(1)} w_j^{(2)} \right] x_i^{(1)} x_j^{(2)} = \sigma^2 x^{(1)^T} x^{(2)}$$

# Linear layer
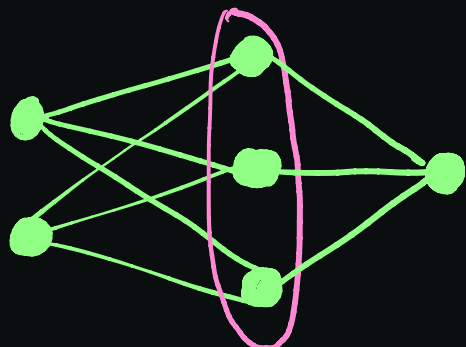
weights $w$

$y_1 = w^T x$

input
$x$

$y_2 = u^T x$

weights $u$

Provided the two weight vectors $w$ and $u$ are sampled independently of each other, then two neurons in a linear layer are just independent linear neurons.

The linear neuron yields a GP $\Rightarrow$ the linear layer does too.

# One hidden layer

nonlinearity

$$y(x) = \sum_{i=1}^{d_1} u_i \, \phi\left(w_i^T x\right)$$

$d_1$ hidden units

second layer weights

weight vector of ith hidden neuron

For inputs $x^{(1)}, \dots, x^{(n)}$, consider vector

$$\left[y(x^{(1)}), \dots, y(x^{(n)})\right] = \sum_{i=1}^{d_1} \left[u_i \phi(w_i^T x^{(1)}), \dots, u_i \phi(w_i^T x^{(n)})\right]$$

If all the weights are drawn iid with finite variance (and $\phi$ does not blow up variances) then each summand is an iid random vector with finite covariance.

So, as $d_1 \to \infty$, $\left[y(x^{(1)}), \dots, y(x^{(n)})\right]$ is Gaussian by the MV-CLT.  21

# **Many hidden layers**



$d_1$ units   $d_2$ units   $d_{L-1}$ units

$x \rightarrow$   $y(x) \rightarrow$

$h_1(x)$   $h_2(x)$   $h_{L-1}(x)$

Again, for inputs $x^{(1)}, x^{(2)}, \ldots, x^{(n)}$, consider vector

$$\left[ y(x^{(1)}), \ldots, y(x^{(n)}) \right] = \sum_{i=1}^{d_{L-1}} \left[ u_i \, \phi\left( h_i^{L-1}(x^{(1)}) \right), \ldots, u_i \, \phi\left( h_i^{L-1}(x^{(n)}) \right) \right]$$

We'd like to take $d_{L-1} \rightarrow \infty$ and apply the MV-CLT as we did for one hidden layer, but first we need to check that for a fixed input $x$, the hidden units $h_1^{L-1}(x), h_2^{L-1}(x), \ldots, h_{d_{L-1}}^{L-1}(x)$ are independent.

Surprisingly, this does hold in the limit that $d_1, d_2, \ldots, d_{L-2} \rightarrow \infty$.

The proof is by induction, using the MV-CLT.

22

# Relu networks

# Relu MLP covariance

- $L$-layer MLP, hidden layer width $\to \infty$
- nonlinearity $\phi(z) := \sqrt{2} \cdot \max(0, z)$
- inputs $x_1, x_2, ..., x_n \in \mathbb{R}^d$ with $\|x_i\|_2 = \sqrt{d}$
- iid Gaussian weights $\mathcal{N}\left(0, \frac{1}{\text{fan in}}\right)$

Define $h(t) := \frac{1}{\pi}\left[\sqrt{1-t^2} + t(\pi - \arccos t)\right]$.

Then for two inputs $x, x' \in \mathbb{R}^d$,

$$\mathbb{E}\left[f(x) f(x')\right] = \underbrace{h \circ ... \circ h}_{L-1 \text{ times}}\left(\frac{x^T x'}{d}\right)$$

the compositional arccosine kernel.

Interpretation: $\frac{x^T x'}{d}$ is the covariance for a linear neuron

The additional layers modify this covariance via the fn. $h$.

## Summary

We found that for a finite collection of inputs, the outputs of a sufficiently wide NN are jointly Gaussian w.r.t. random sampling of the weights.

The covariance matrix of this Gaussian depends on how the network architecture transforms the input correlations.

# Next lecture

The first application of our tools:

— optimisation of neural networks.