

Shift gears a bit and study results that are more directly applicable to classification

First paper: Montúfar, Pascanu, Cho, Bengio, 2014
"On the number of linear regions in deep neural networks"

See also: Pascanu, Montúfar, Bengio, 2013
"On the # of response regions of deep feed forward networks w/ piecewise linear activations"
and more recent papers!

Idea: Count the # of piecewise linear parts of a function $f(x; \theta)$ that a ReLU network can implement

Summary of result: Deep ReLU networks can implement functions $f(x; \theta)$ with exponentially more piecewise linear regions than shallow networks w/ the same # of hidden neurons.

Now show Figure 1 from the paper

Types of networks we will study: $f(\cdot; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^{\text{out}}$
(i) Shallow, one layer:

$$f(x; \theta) = f_{\text{out}} \circ \sigma \circ A(x), \quad A(x) = Wx + b$$

where: $W \in \mathbb{R}^{d_1 \times d}$, i.e., d_1 neurons
 $b \in \mathbb{R}^{d_1}$

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

$$f_{\text{out}} : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{\text{out}}, \text{ e.g.}$$

- Regression, $f_{\text{out}} : \mathbb{R}^{d_1} \rightarrow \mathbb{R}$ and $f_{\text{out}}(z) = \langle z, \alpha \rangle$
- Classification, $f_{\text{out}} : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{\# \text{ classes}}$ and $f_{\text{out}} = \text{softmax}$

(ii) Deep, L layers:

$$f(x; \theta) = f_{\text{out}} \circ \sigma \circ A_L \circ \dots \circ \sigma \circ A_1(x)$$

where

$$A_\ell(x) = W_\ell x + b_\ell, \quad W_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}, \quad d_0 = d$$
$$b_\ell \in \mathbb{R}^{d_\ell}$$

$$f_{\text{out}} : \mathbb{R}^{d_L} \rightarrow \mathbb{R}^{\text{out}}$$

Rest same as shallow network.

Def: Let $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a piecewise linear function. A linear region of f is a maximal connected subset of the input space \mathbb{R}^d on which f is linear.

Note: For ReLU networks $f(x; \theta)$, the dimension of each linear region is d .

Results for shallow networks

Let $f(x; \theta)$ be a shallow ^{ReLU} network w/ one hidden layer and d_1 neurons. Then the maximum # of linear regions of $f(x; \theta)$ is

$$\text{max \# of linear regions} = \sum_{i=0}^d \binom{d_1}{i}$$

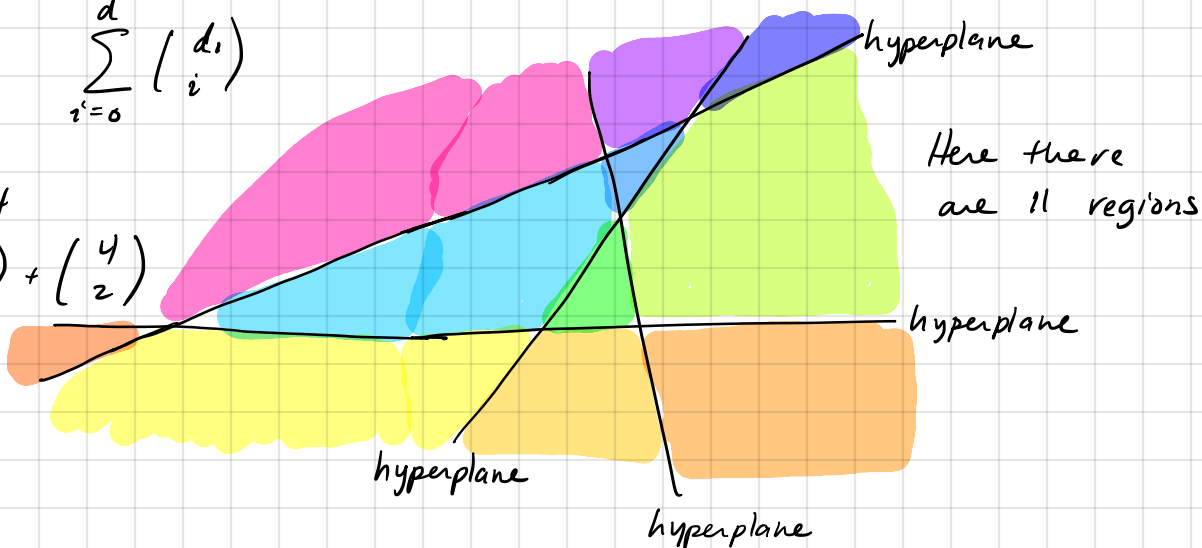
Idea of proof: Suppose you have d_1 hyperplanes in \mathbb{R}^d . The maximum # of regions these hyperplanes can divide \mathbb{R}^d into is

$$\sum_{i=0}^d \binom{d_1}{i}$$

In this example:

$$d=2, d_1=4$$

$$\begin{aligned} \text{max \#} &= \binom{4}{0} + \binom{4}{1} + \binom{4}{2} \\ &= 1 + 4 + 6 \\ &= 11 \end{aligned}$$



Results for deep networks

Main result: Let $f(x; \theta)$ be deep ReLU network w/ L hidden layers and $d_l \geq d$, $1 \leq l \leq L$, neurons in each layer. Then the maximal # of linear regions in $f(x; \theta)$ is at least:

$$\text{max \# of linear regions} \geq \left(\prod_{l=1}^{L-1} \left\lfloor \frac{d_l}{d} \right\rfloor^d \right) \sum_{i=0}^d \binom{d_L}{i} \quad (*)$$

Come back to this at the end

This part just follows from one layer analysis

Asymptotic comparison of shallow to deep networks

Let the input dimension d be fixed.

(i) For the shallow network, give it $d_1 = Ln$ neurons

(ii) For the deep network w/ L layers, let each layer have $d_l = n$ neurons

Therefore both networks have the same # of neurons.

Using the shallow network result, we see:

$$\max_{\text{linear}} \# \text{ regions shallow network} = O(L^d n^d) \leftarrow \begin{array}{l} \text{bounded from above} \\ \text{asymptotically} \\ \text{as } L, n \rightarrow \infty \end{array}$$

polynomial in n
polynomial in L

Using the deep network result, we see:

$$\max_{\text{linear}} \# \text{ regions deep network} = \Omega\left(\left(\frac{n}{d}\right)^{(L-1)d} n^d\right) \leftarrow \begin{array}{l} \text{bounded from} \\ \text{below asymptotically} \\ \text{as } L, n \rightarrow \infty \end{array}$$

polynomial in n (same as shallow)
exponential in L (larger than shallow!)

Back to (*)

$$(*) : \max \# \text{ linear regions for deep network} \geq \underbrace{\left(\prod_{l=1}^{L-1} \left\lfloor \frac{d_l}{d} \right\rfloor^d \right)}_{(**)} \sum_{i=0}^d \binom{d}{i}$$

Idea behind (**):

(1) Each layer can divide a single linear region into

$\left\lfloor \frac{d_l}{d} \right\rfloor^d$ linear regions

(2) So the first layer can create $\left\lfloor \frac{d_1}{d} \right\rfloor^d$ linear regions

(3) The 2nd layer can create $\left\lfloor \frac{d_2}{d} \right\rfloor^d$ linear regions from

each of the linear regions created by the 1st layer.

Thus the total # of linear regions is:

$$\left\lfloor \frac{d_1}{d} \right\rfloor^d \cdot \left\lfloor \frac{d_2}{d} \right\rfloor^d$$

(4) Continue for the first $L-1$ layers, you get (**).