

Exponential Gains from Depth

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Representation Learning

1. Greedy Layer-Wise Unsupervised Pretraining
2. Transfer Learning and Domain Adaptation
3. Semi-supervised Disentangling of Causal Factors
4. Distributed Representation
- 5. Exponential Gains from depth**
6. Providing Clues to Discover Underlying Causes

Topics in Exponential Gains in Depth

- Efficiency of multilayer perceptrons
- Expressive power of deep networks
- Minimum depth required
- Expressive power of CNNs

Theoretical Underpinnings

- Multilayer perceptrons are universal approximators
 - Can approximate most functions given enough hidden units upto any non-zero tolerance
- There exist families of functions which can be approximated efficiently with depth $> d$, but which require a much larger model if depth $\leq d$
 - In many cases, no. of hidden units required by the shallow model is exponential in n

Deep is more efficient than Shallow

- Result for MLPs

- Functions are represented by smaller deep networks compared to shallow networks

- No of linear regions with d inputs, depth l and n units per hidden layer is exponential in l

$$O\left(\binom{n}{d}^{d(l-1)} n^d\right)$$

- Decrease in model size leads to statistical efficiency
 - This section looks at how similar results apply to other distributed hidden representations

Size of h for Universal approximators

- Composition of nonlinearities and hierarchy of reused features can give exponential boost to statistical efficiency, on top of exponential boost given by a distributed representation
 - Many networks (e.g., with saturating nonlinearities, Boolean gates, sum/products, or RBF units) with a single hidden layer are universal approximators
 - A model family that is a universal approximator can approximate a large class of functions (including all continuous functions) up to any non-zero tolerance level, given enough hidden units.
 - However, required no of hidden units h is large

Expressive power of deep networks

- There are families of functions that can be:
 1. Efficiently represented by architecture of depth k
 2. But would require an exponential number of hidden units (wrt input size) with insufficient depth (depth 2 or depth $k-1$)

Minimum depth required

- Theoretical result based on using a Sum-Product Network (SPN)
 - Which use polynomial circuits
 - Models compute a probability distribution over a set of random variables
- Showed that there exist probability distributions for which a minimum depth of SPN is required to avoid needing an exponentially large model
 - Later showed that there are significant differences between every two finite depths of SPN
 - Some constraints used to make SPNs tractable may limit their representational power

Expressive power of CNNs

- Theoretical result based on families of deep circuits related to convolutional networks
- There is an exponential advantage for the deep circuit even when the shallow circuit is allowed to only approximate the function computed by the deep circuit