

Andrew Saxe

Theory of Rep Learning in Complex Architectures

What do DL nets actually learn?

Semantic meaning # topic

Deep linear Net

— gains nothing from depth

SVD change

— cross assoc matrix of
items and their properties

show to V

— Decomposes the matrix
into

shit. I didn't get
the nomenclature.

D : size of net
(layers)

$l = 1 \dots D$

w = weights?

Exact solutions

— about lost here

Hidden 1D Layers

— Exposing hidden linear mechanisms is a DLN

like ^{lying} an ocean
feeling the waves push
through your body
also made of water

Error surface

— Getting stuck on initial plateaus

— DLNs have plateaus

SVD change of Variables

~~scribble~~

Mostly been talking about speed & scaling of training speed

~~scribble~~

Semantics & classification

Semantic development

I need to read
Chomsky

D linear net model

- good graph for V

~~scribble~~

Diverse Structures

- Nice illustration

~~scribble~~

Prog Diff

- Theory children learn broad categories before finer distinctions

~~~~~

Idiosyncratic Individuals

- Is he arguing for an "ego"?

- or hidden layer as "ego"?

~~scribble~~

Transformer Heads

# Neural Race Reduction

## Mesoscale Architecture

- "We don't have good ideas of how architecture represents" !

Mupkes et al. 2020

# read

## Theoretical Analysis

- I think this maybe a job for me
- Inference of network activity?

## Gated D Linear Net

- adding scalar gating variables "g"
- multiply by gating var @ weight & @ node

## Gating input

- adding gating variables with inputs

## Grad. Des.

- Compressed path notation
- to describe "path": the full

- Pathway counting logic

→ More shared pathways = more learning

# Neural Race Reduction

- I think we are reducing a <sup>ReLU</sup> network to its effective parts via studying effective paths

## Lazy v. Feature learning

- "Large networks w/ large weights learn better" ?

- lazy is not as useful, more static

## Routing net

I don't get how these are not oversaturated.

like, yes we want as much multiplexing/  
generalizability/  
overlap of calculations  
but ...

How can they do this reduction easily?

I guess they mentioned it with the node contrib feature stuff

- gating variables tell you the importance of edges

Saxe et al. t6p, #read



