

Neural ODE Problem Setup

- We have a func $F(x, t)$
we want to train a dynamic NN to compute $F(x, t_i)$ given $F(x, t_0)$ and samples at $F(x, t)$
- an initial value problem
 $\frac{dF}{dt} = f = F(x, t, \phi)$ derivative of F w.r.t t
whatever time step
can be solved by ODE solver
But we don't know f so we use a NN to learn it
 T analytical

Deep Learning Theory and Applications

Representation Learning in Neural Networks

Yale

CPSC/AMTH 663

Summary How to design NN

- ① initial: deep and wide
 - ② prune (neuron or edge)
 - ③ narrower subsequent hidden layers : JL Lemma
- } lottery Ticket hypothesis





What is happening inside a neural network?

- The output of a neural network is a function
- We learned last time that it can be ANY function
 - Should be able to approximate it
- Now what about the other layers?
- How are they facilitating the learning?



Dimensionality Reduction

- Neural networks typically learn a dimensionality reduced embedding in their internal layers
- These embeddings form GOOD projections of the data
 - Projections that maintain distances between datapoints
 - Projections that can show groupings, clusters etc $\|x_1 - x_2\| \approx \|f(x_1) - f(x_2)\|$
- This can be shown mathematically

keep important info and remove useless info

Johnson Lindenstrauss Lemma 定理



sample size



Theorem 2.1. For any $0 < \epsilon < 1$ and any integer n , let k be a positive integer such that

$$k \geq 4(\epsilon^2/2 - \epsilon^3/3)^{-1} \ln n. \quad (2.1)$$

Then for any set V of n points in \mathbf{R}^d , there is a map $f: \mathbf{R}^d \rightarrow \mathbf{R}^k$ such that for all $u, v \in V$,

$$(1 - \epsilon)\|u - v\|^2 \leq \|f(u) - f(v)\|^2 \leq (1 + \epsilon)\|u - v\|^2.$$

squared distance

close to original distance

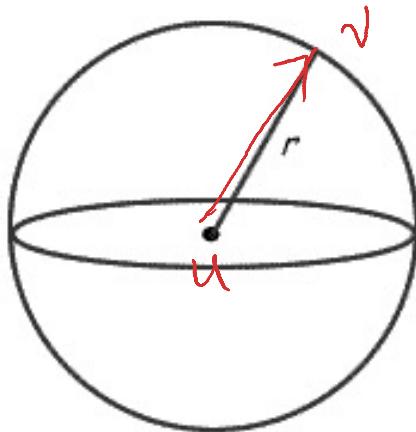
Furthermore, this map can be found in randomized polynomial time. $T(O^n)$

d : input size

k : embedding size



Length of projected unit vector



- What is the expected length of a unit vector in d dimensions when it is projected down to K ?
- If several such vectors are projected down would their relative magnitudes be preserved?



Proof Concept

- It shows that the squared length of a random vector is sharply **concentrated around its mean** when the vector is projected onto a random k-dimensional subspace. Specifically, with probability $O(1/n^2)$, its (scaled) length is not distorted by more than $(1 + \epsilon)$
- Let X_1, \dots, X_d be d independent Gaussian $N(0, 1)$ random variables, and let $Y = \frac{1}{\|x\|} (X_1, \dots, X_d)$. Y is a point chosen uniformly at random from the surface of the d -dimensional sphere S^d . Let the vector $Z \in \mathbb{R}^k$ be the projection of Y onto its **first k coordinates**, and let $L = \|Z\|^2$.
- The expected squared length of Z is $\mu = E[L] = k/d$
 - And the probability is tightly concentrated around the mean



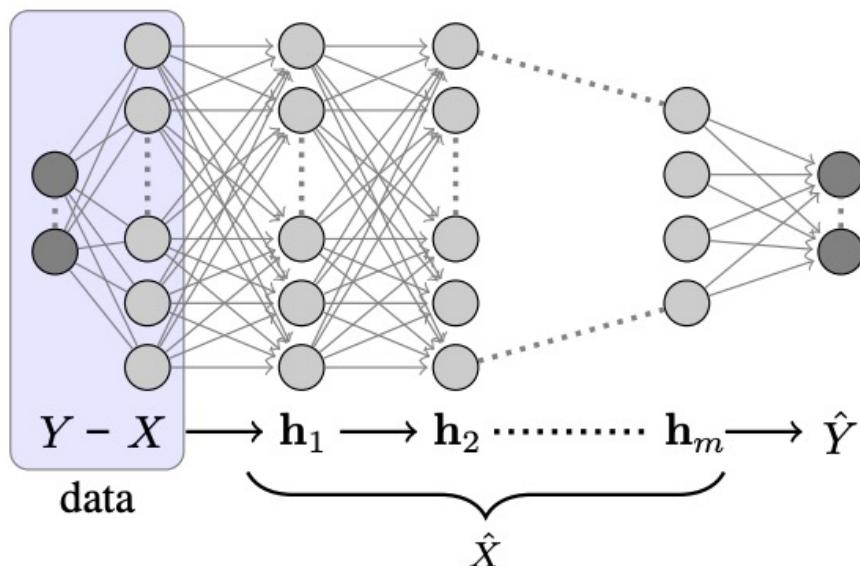
Johnson-Lindenstrauss lemma

- How does this connect to neural networks?
- This means that **random initialization** helps find good results, local minima are pretty good
- High probability of cluster patterns, manifolds, existing *bc preserve distance will preserve most info*
- SGD will only improve upon it.



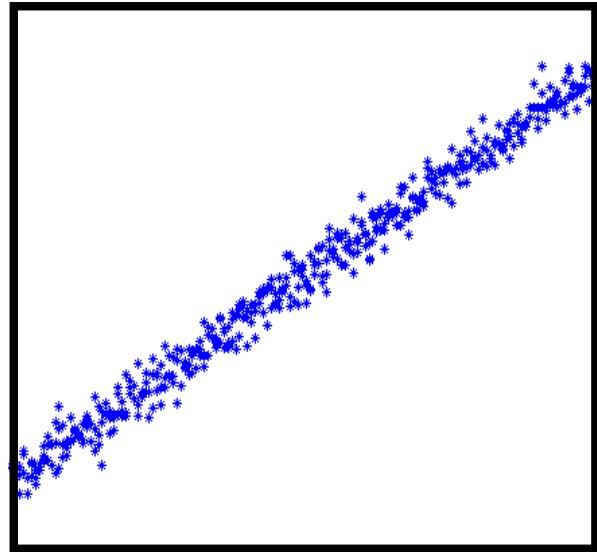
Deep Projections

- We have projections, but why make them deep?
- What do subsequent layers do?

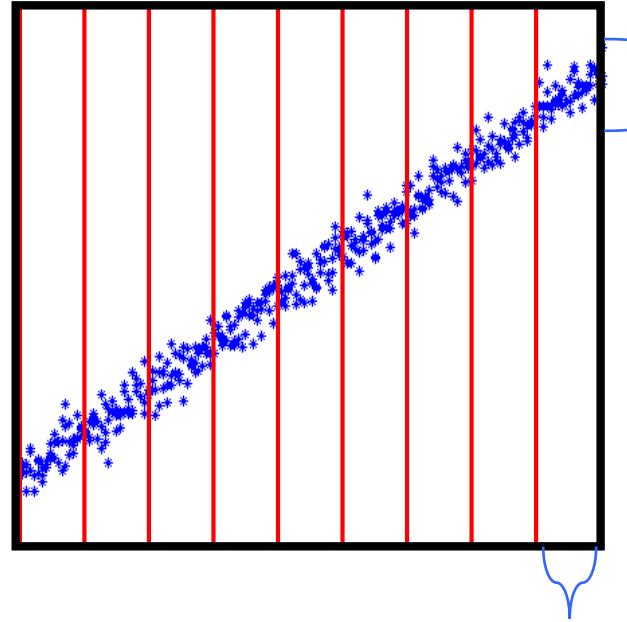




Entropy as a measure of uncertainty



Y Spread



Y Spread
in X-Slice

Measure of Uncertainty in a random variable.

$$-\sum_{i=1}^n P(x_i) \log_b P(x_i),$$

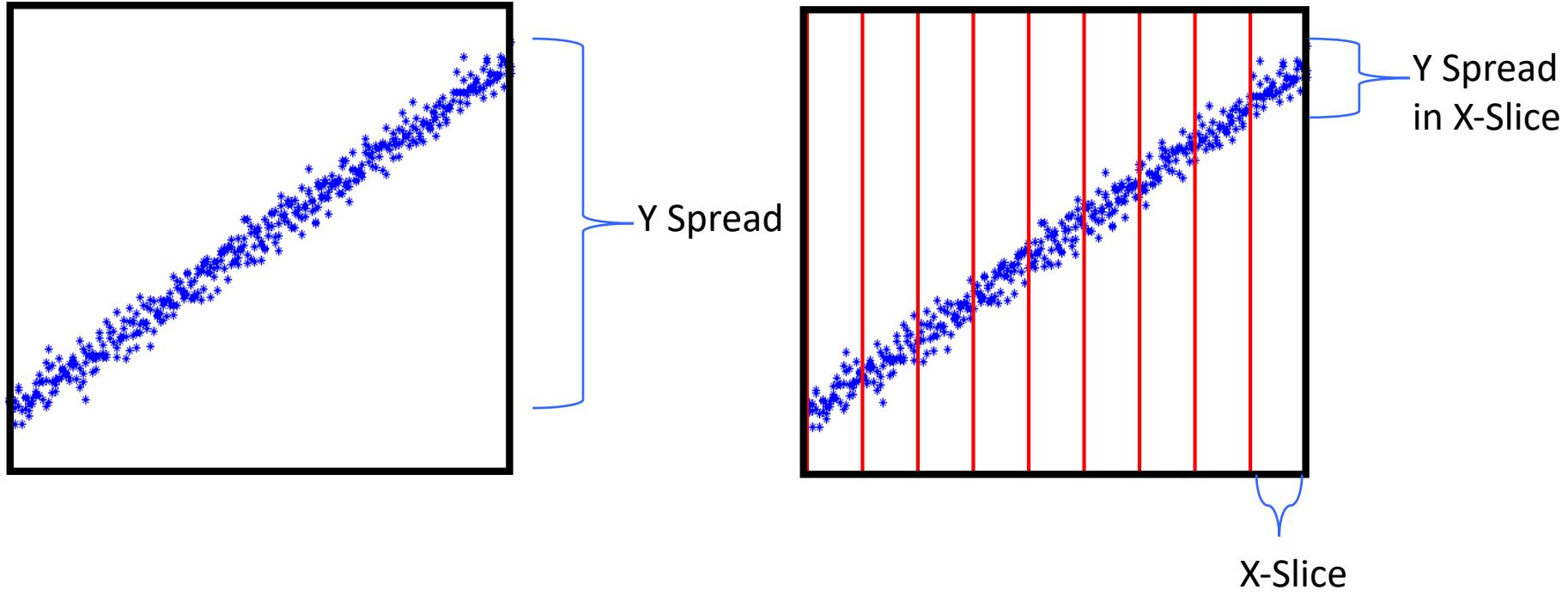
Shannon Entropy

X-Slice

Units of bits tells us how many bits are needed to represent the outcome



Mutual Information



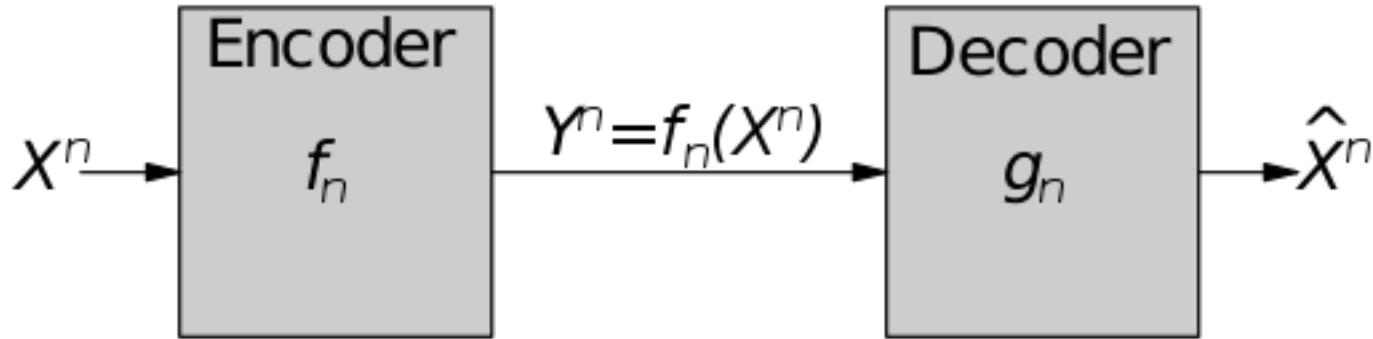
How much info X on Y

$$\text{Mutual Information: } I(X,Y) = H(Y) - H(Y|X)$$

↑
unconditional entropy ↑
conditional entropy



Rate Distortion Theory



$$\inf_{Q_{Y|X}(y|x)} I_Q(Y; X) \text{ subject to } D_Q \leq D^*.$$



Information Bottleneck Principle

- Introduced by Tishby et al. uses information theory to understand what about the ~~input~~ is relevant to the output ^X _Y
- Given a joint distribution $p(X,Y)$, the ~~information relevant to Y within X~~ ^{or internal layers to output} is given by $I(X,Y)$ the *mutual information* between X and Y
- Therefore, a good compression of X should capture the ~~mutual information~~ and discard the remaining information $X - I(X,Y)$



Information Bottleneck Loss

- Suppose the relevant information about Y is captured by minimum sufficient statistics \hat{X} (compressed repr of X)
- We can define a loss information bottleneck loss

$$\mathcal{L}[p(\hat{x}|x)] = I(X; \hat{X}) - \beta I(\hat{X}; Y)$$

high loss: $\begin{cases} \text{large } I(X; \hat{X}) & \hat{X} \approx X \\ \text{small } I(\hat{X}; Y) & \hat{X} \text{ has no info on } Y \end{cases}$ compression rate info preserv

- This defines a tradeoff finding a compressed representation of X , and retaining information about Y based on the parameter Beta
- Compression rate vs information preservation

$$R^{(I)}(I_Y) = \min_{p(\hat{x}|x): I(\hat{X}; Y) \geq I_Y} I(X; \hat{X})$$

threshold

minimize $I(X; \hat{X})$ while maintaining $I(\hat{X}; Y) \geq$ some threshold I_Y



IB Loss

- The information bottleneck loss can be written as (only differ by constant)

$$\tilde{\mathcal{L}}[p(\hat{x}|x)] = I(X; \hat{X}) + \beta I(X; Y|\hat{X}) \quad \text{condition on } \hat{X}$$

- With this you can define a distortion

$$D_{IB} = E[d_{IB}(X, \hat{X})] = I(X; Y|\hat{X})$$

- Distortion is the expected Kullback-Liebler divergence

$$d_{IB}(x, \hat{x}) = D[p(y|x) \| p(y|\hat{x})]$$



DNNs as Markov Chains

- DNNs use multiple internal layers to compress representations
- These can be thought of Markov Chains and reasoning about them can be done via the Data Processing Inequality

$$I(Y; X) \geq I(Y; \mathbf{h}_j) \geq I(Y; \mathbf{h}_i) \geq I(Y; \hat{Y})$$

input hidden layer $j < i$ predicted output

- Achieving equality only happens when maximizing information to the output layer but also minimizing as much as possible information between subsequent layers

$$I(\mathbf{h}_{i-1}; \mathbf{h}_i)$$



Relevant quantities

- It has been shown that can bound classification error in tasks with multiple classes $I(Y; \mathbf{h}_i)$
depends on info
- This is the mutual information between the representation in layer i and the class labels in Y
- This is the minimal description length $\# \text{ bits}$
 $I(\mathbf{h}_{i-1}; \mathbf{h}_i)$ *Info between hidden layers*
- Each layer can be compared to the theoretical optimal layer for the loss $I(\mathbf{h}_{i-1}; \mathbf{h}_i) + \beta I(Y; \mathbf{h}_{i-1} | \mathbf{h}_i)$



Optimal Solutions for IB Loss

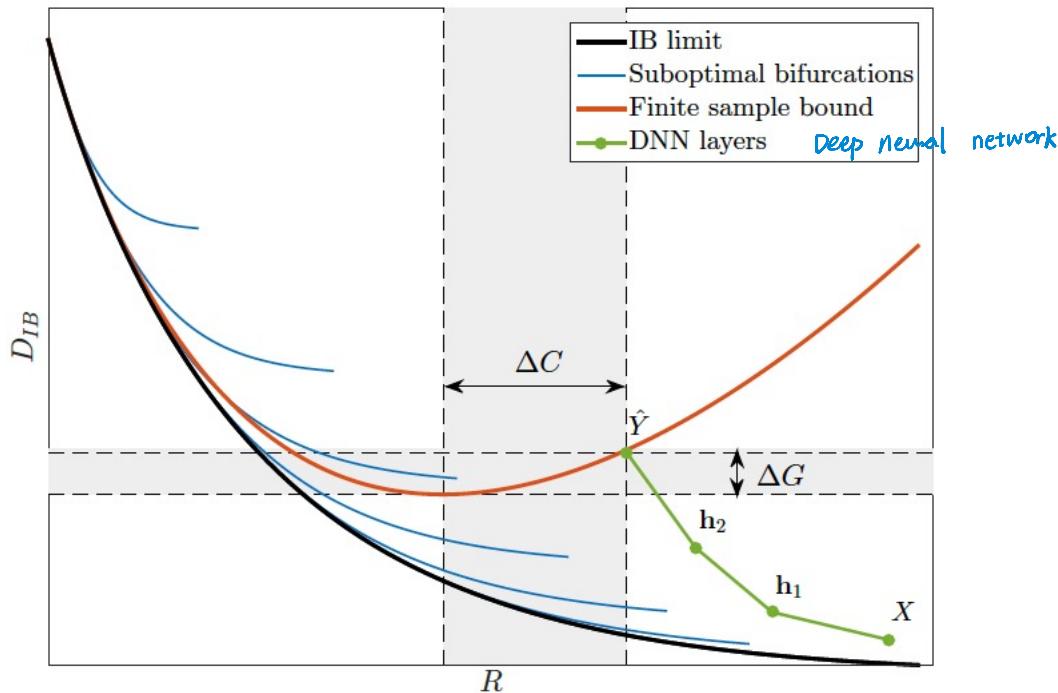
- Can be found by iterating these equations

$$\begin{aligned} p(\hat{x}|x) &= \frac{p(\hat{x})}{Z(x;\beta)} \exp(-\beta D[p(y|x) \| p(y|\hat{x})]) \\ p(y|\hat{x}) &= \sum_x p(y|x) p(x|\hat{x}) \\ p(\hat{x}) &= \sum_x p(x) p(\hat{x}|x) \end{aligned}$$

- This gives rise to an optimal rate distortion curve (not one that is reached by DNNs)



Rate Distortion Curve



β is the negative slope of this curve

交叉

a big change in info

Bifurcations happened where there are topological changes in the representation
Green are DNN layers are generally higher than the optimal



Multiple DNN Layers

- Each layer can only increase the distortion (based on data processing inequality)----can't add lost information
- However layers can compress information better in order to meet the finite sample bound curve
- This is the bound you can get using only a finite sample

$$I(X; \hat{X}) \leq \hat{I}(X; \hat{X}) + O\left(\frac{K}{\sqrt{n}}\right) \quad K = |\hat{\mathcal{X}}|$$



IB-based Reason for Multiple Layers

- Each layer in a neural network effectively only linearly separates the information in the previous layer
- Recall: One layer neural networks are equivalent to SVMs
- This is because layers are non-linear functions of dot products
- Linear separability is only possible when input units are conditionally independent given output classification
- This is not necessarily true of the data distribution, multiple hidden layers are required!



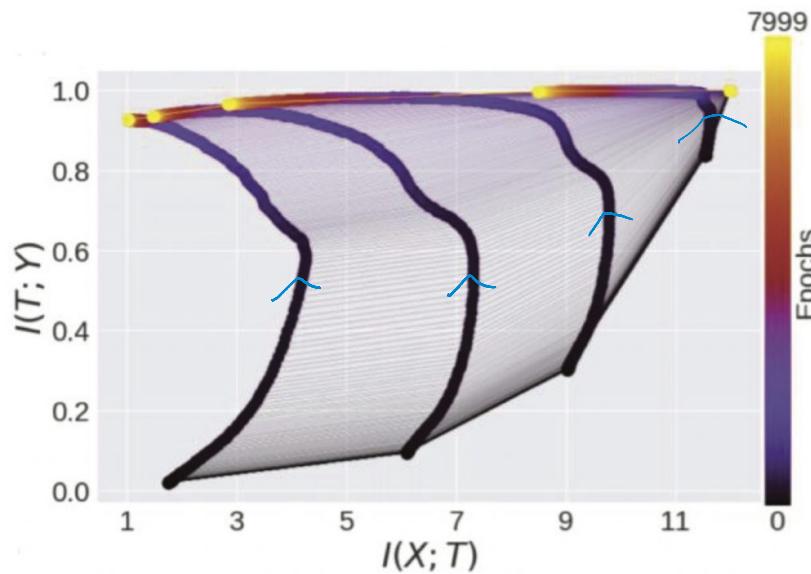
Phase Transitions

- Interesting theory broached by Tisby et al. is that one could find the correct architecture
- Number of layers and number of neurons in each layer by looking to see phase transitions of the rate distortion curve using bifurcation analysis
- Namely finding critical values of β where one can move to simpler representations



During Training

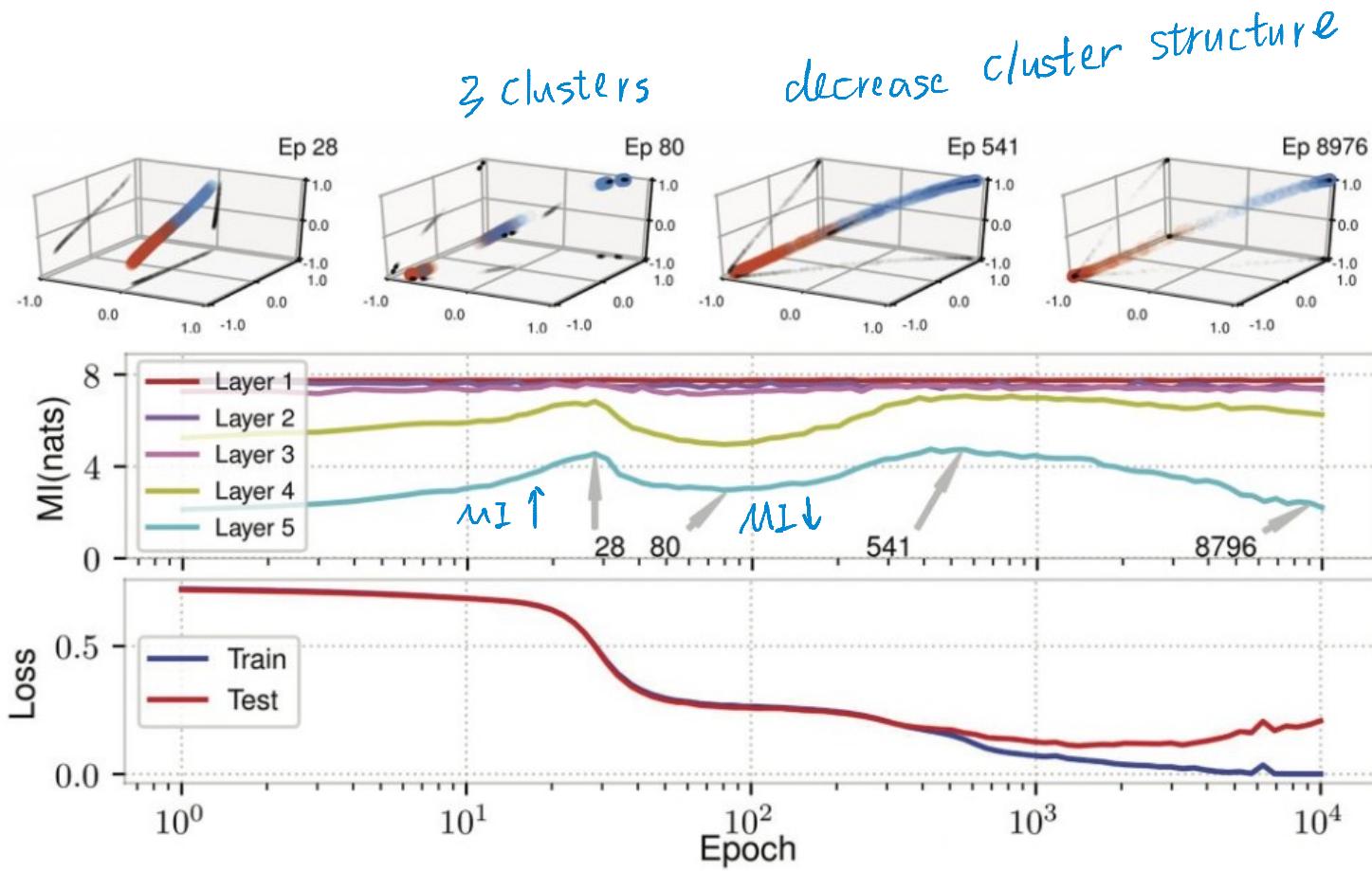
- It has been observed that the mutual information between the input layer and the hidden layers **initially increases** as the neural network learns to encode the input
- Then **falls rapidly** as it learns to compress this information



T: layer T
During training
 $I(X; T)$ 先↑后↓
 $I(T; Y)$ 一直↑



Evolution of representation



Rise and fall of MI corresponds to how clustered the representation will be



How do representations evolve?

- What happens during training
- We have shown that random initializations start well
where do they go after training?

M-PHATE

Visualizing Internal Representations
(Gigante,S. et al. NeurIPS 2019)



Main idea: Multi-slice kernel allows for PHATE to visualize dynamic behaviorin neural networks



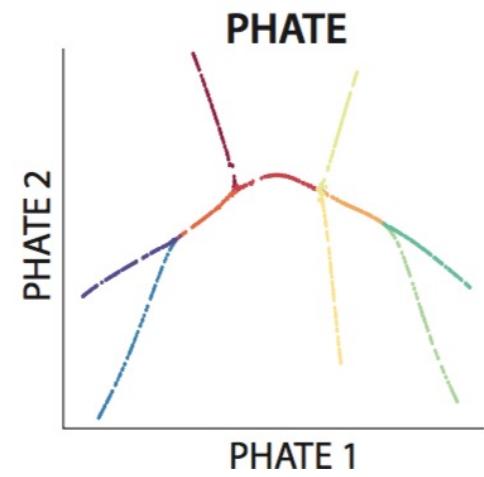
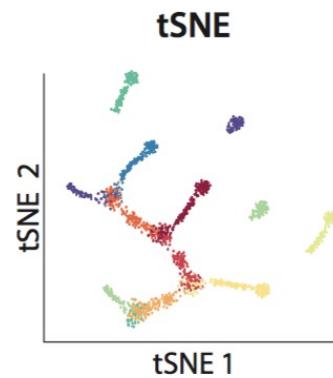
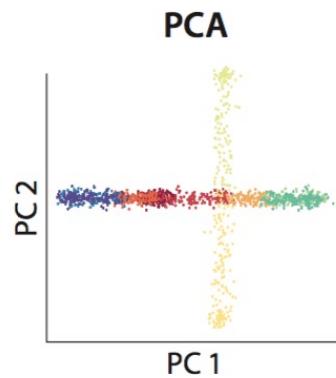
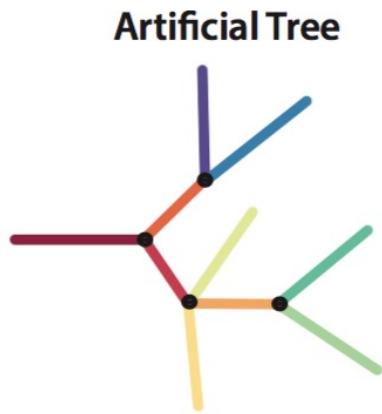
Scott Gigante



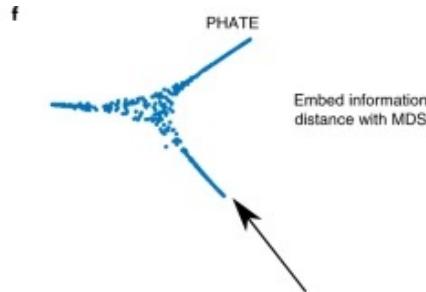
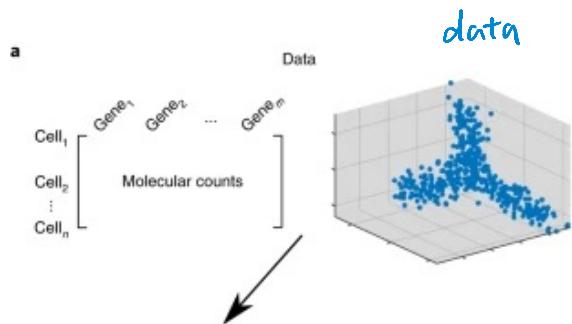
Gal Mishne

Scott Gigante
Adam Charles
Gal Mishne

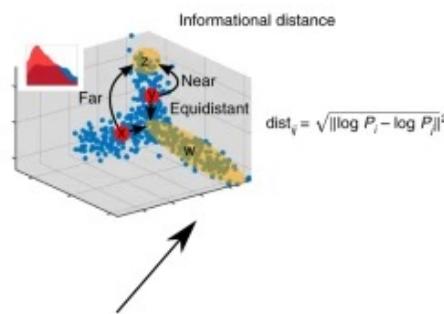
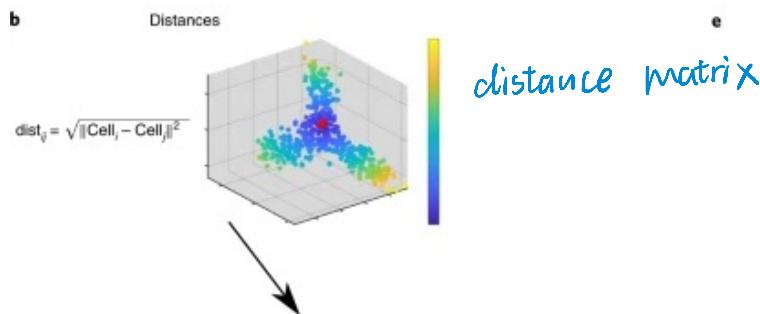
Visualization methods fail to capture global structure in high dimensional data



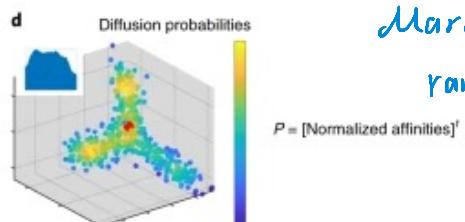
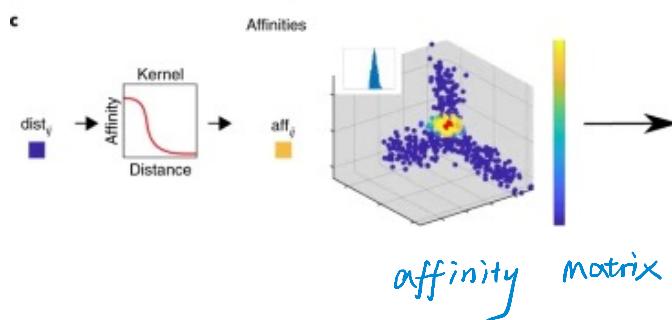
*keep distance in
manifold*



visualize in 2D



M-divergence between
probabilities



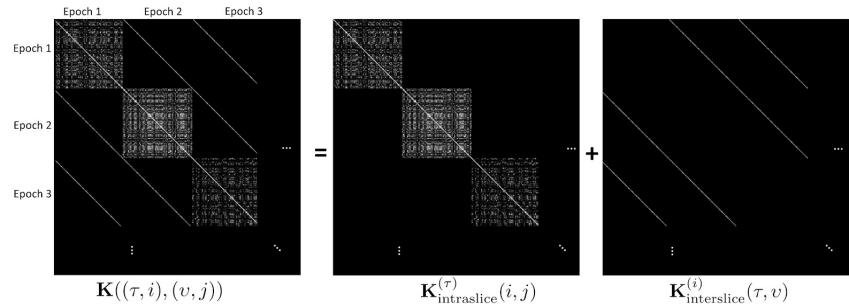
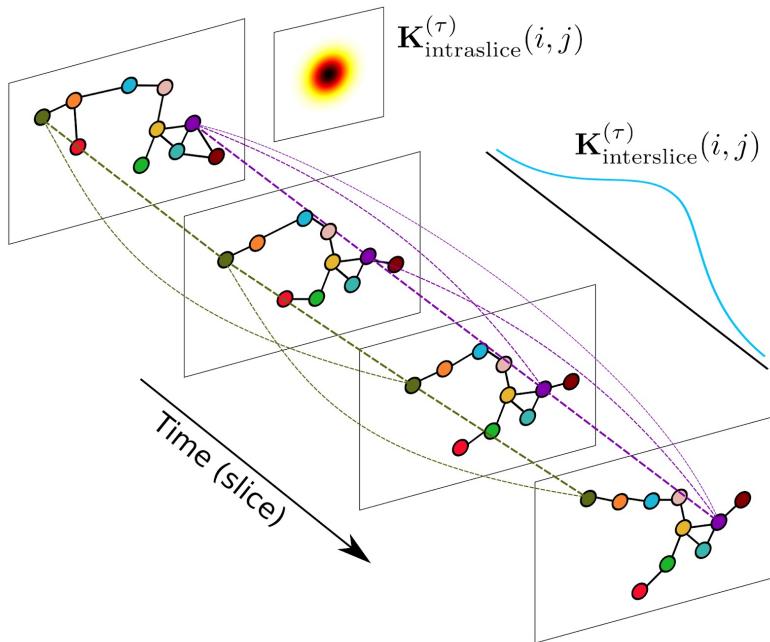
Markov normalized matrix
random walk

M-PHATE Neural Network

Multi-slice PHATE Learning

let PHATE sweep thru time slice

induce affinity between neural itself at different time points



Learning vs Memorization in NNets

How to make NN memorize rather than learn sth

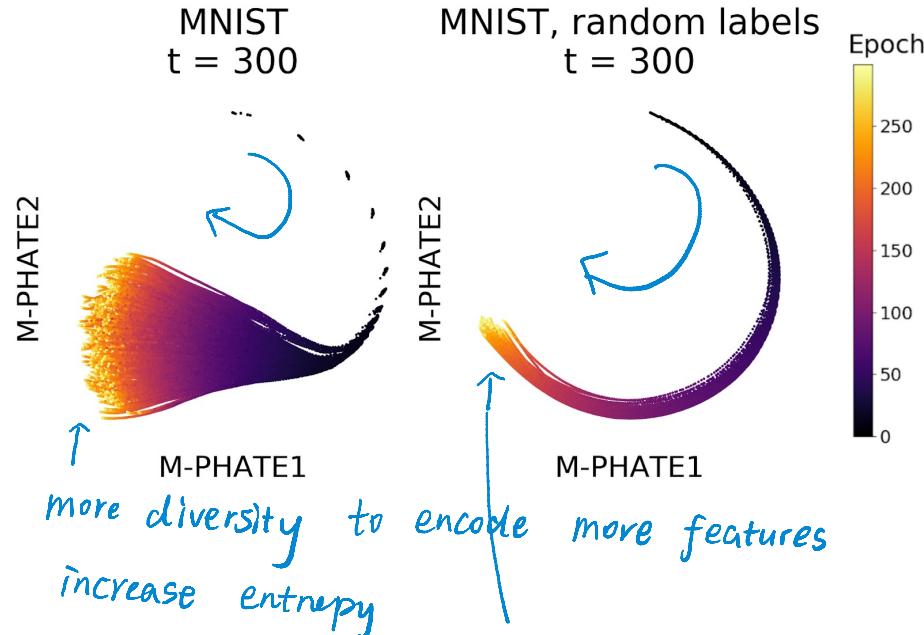
How to make NN memorize rather than learn sth?

A: give scramble labels rather than true labels (meaningful)

无序的

visualize internal layers at epoch

Mnist handwritten digits

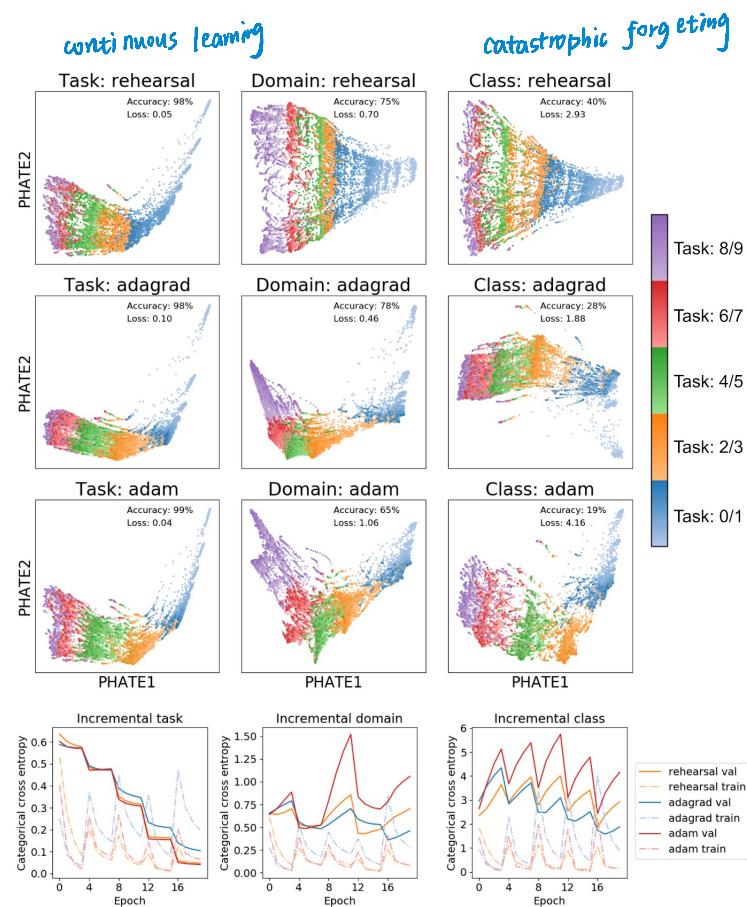


more diversity to encode more features
increase entropy

One neuron is 1 all others 0
bc it just memorize one input

Learning vs. Catastrophic Forgetting

when train NN to do 2nd task
it forget 1st task



Expanding and shrinking

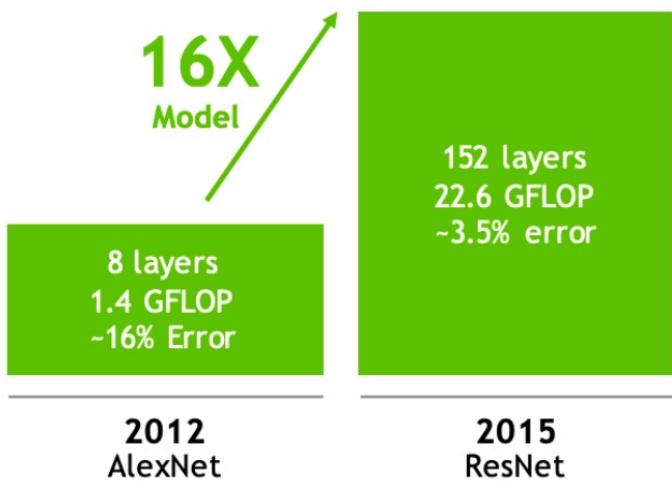
- We have shown that we can use depth to increase the power of a neural network
- Can we also shrink or prune a network?
- There could be many reasons to shrink or prune a network



Reducing Model Size

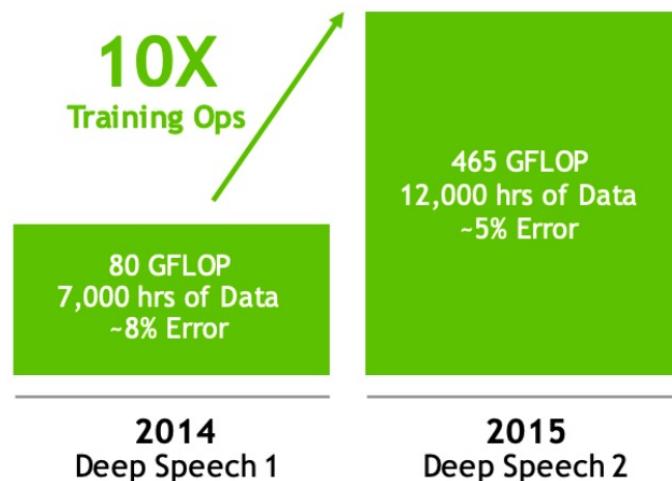
Larger Model => More Memory References => More Energy

IMAGE RECOGNITION



Microsoft

SPEECH RECOGNITION



Baidu



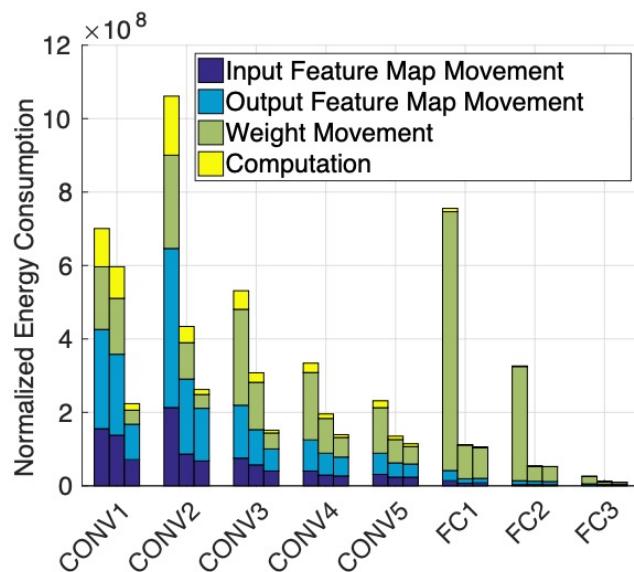
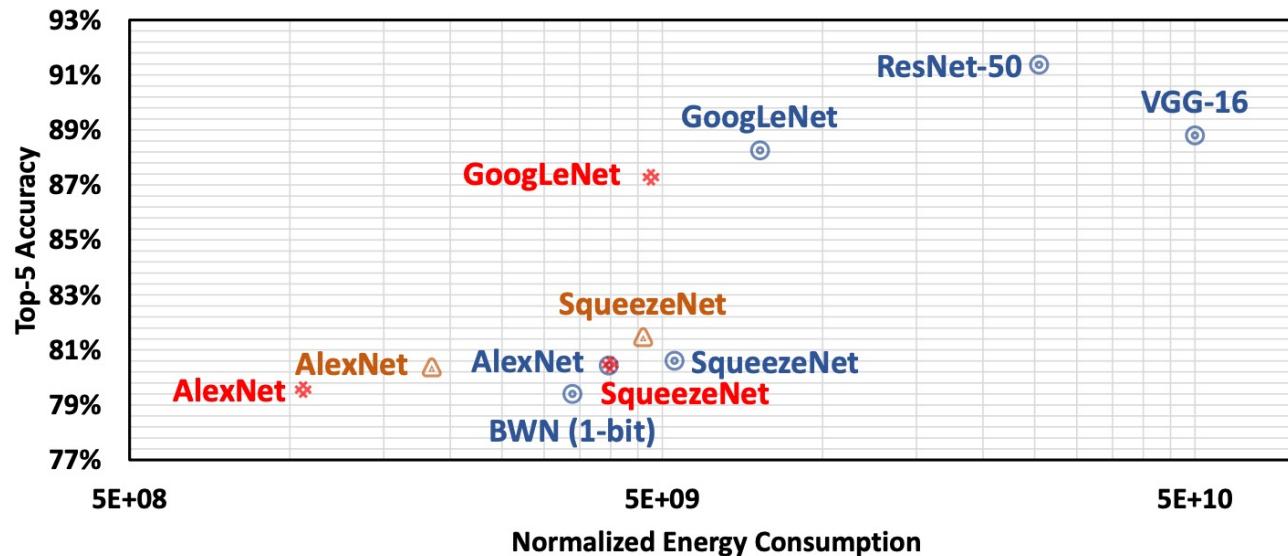
Difficulty Training Large Models

	Error rate	Training time ↑
ResNet18:	10.76%	2.5 days
ResNet50:	7.02%	5 days
ResNet101:	6.21%	1 week
ResNet152:	6.16%	1.5 weeks

Training time benchmarked with fb.resnet.torch using four M40 GPUs



Energy consumption





Biological Inspiration for Pruning



Newborn



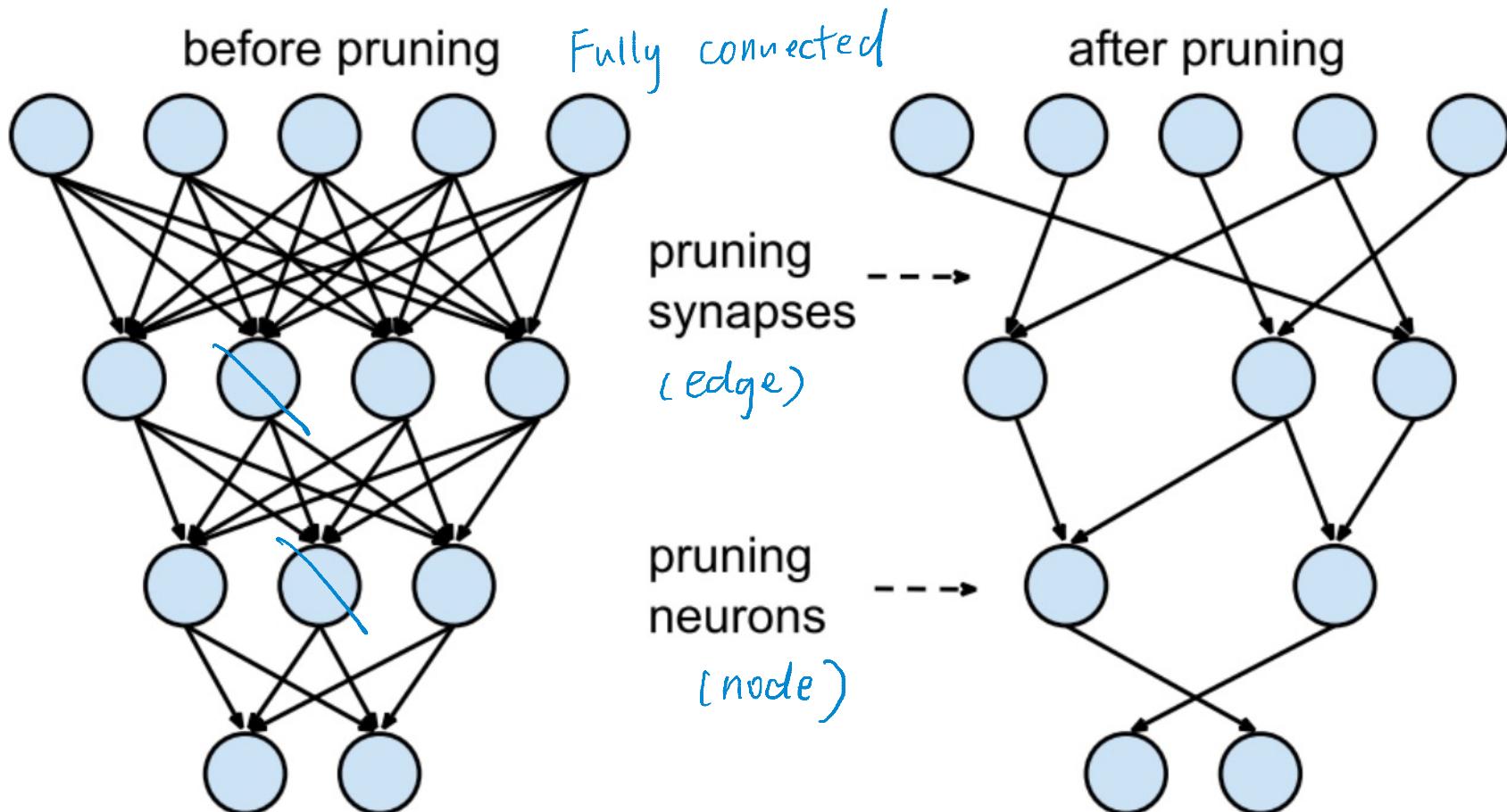
1 year old



Adolescent



Neural Network Pruning





Weight Pruning

- Set individual weights in the weight matrix to zero. This corresponds to deleting connections
- To achieve k% sparsity rank the individual weights in weight matrix W according to their magnitude, and then set the smallest k% to zero.

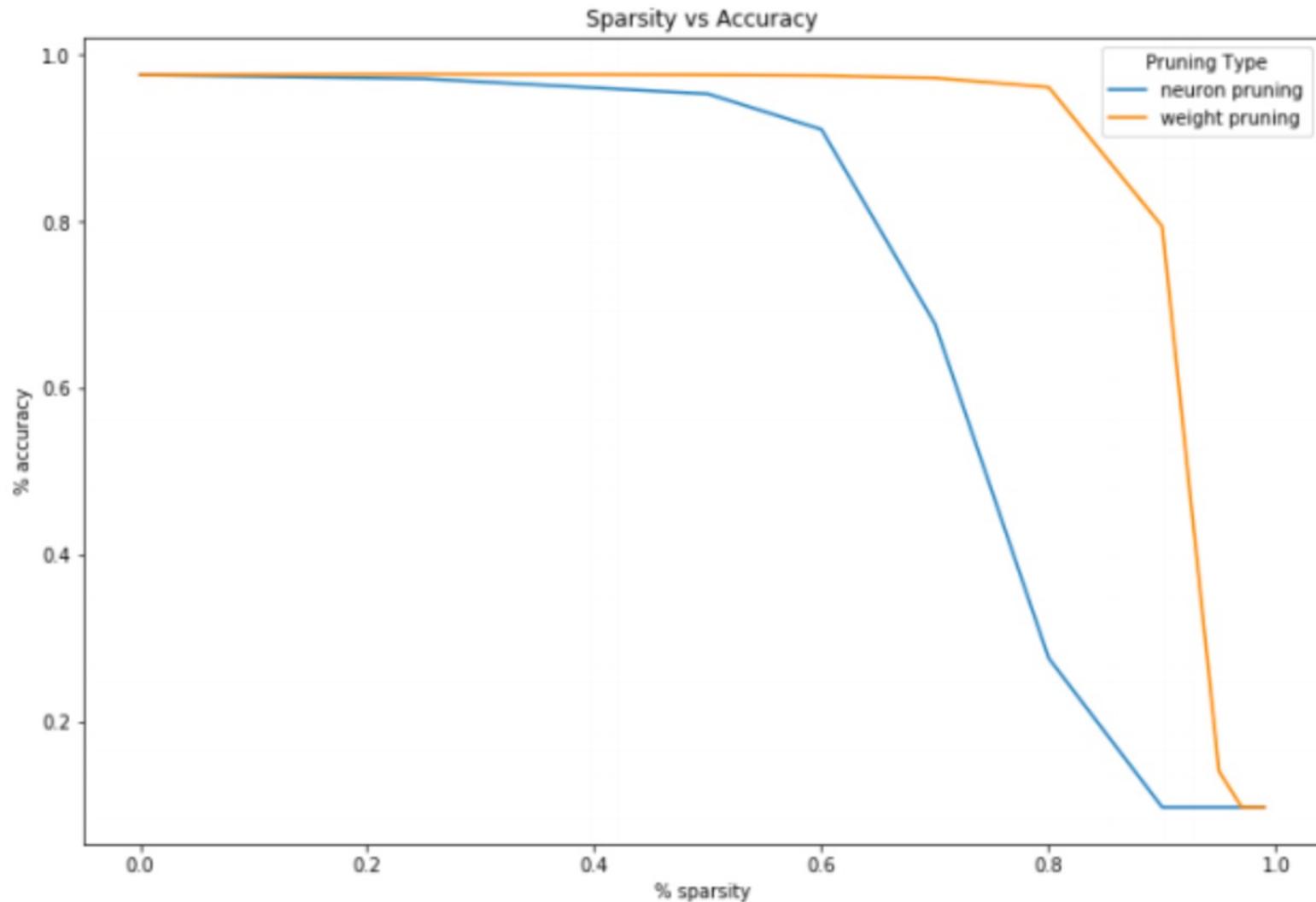


Neuron Pruning

- Set **entire columns** to zero in the weight matrix to zero, in effect deleting the corresponding output neuron.
- Here to achieve **sparsity of k%** we rank the columns of a weight matrix according to their L2-norm and delete the smallest k%.



Pruning often Maintains Accuracy





Lottery Ticket Hypothesis

why?

- Dense, randomly-initialized, feed-forward networks (all tickets) contain subnetworks (winning tickets) that—when trained in isolation—reach test accuracy comparable to the original network in a similar number of iterations.
- The winning tickets we find have won the initialization lottery: their connections have initial weights that make training particularly effective. but if initialize all weights to 0 no such lottery effect
- However, architectures uncovered by pruning are harder to train from the start
 - They reaching lower accuracy than the original networks

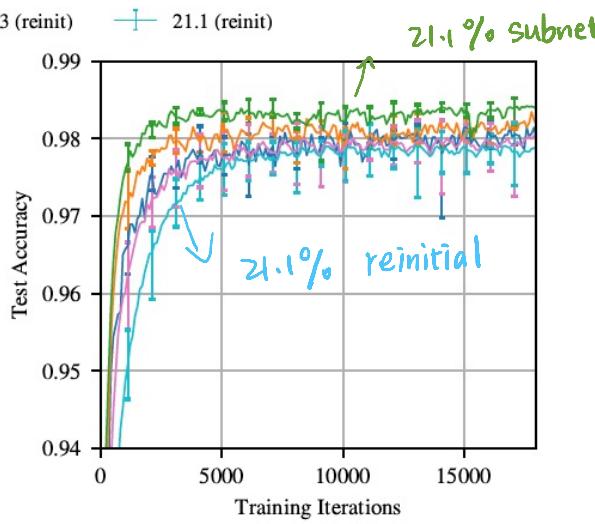
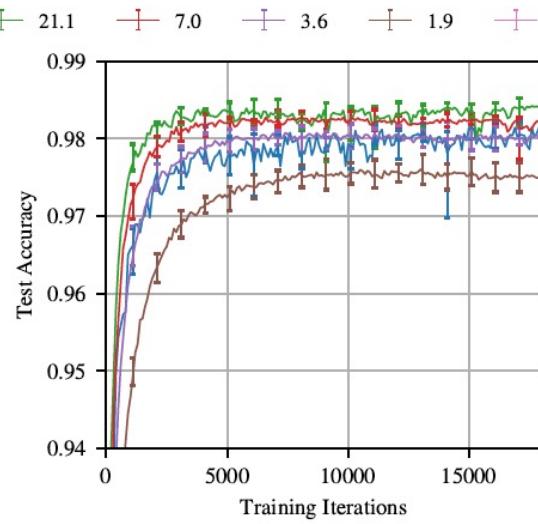
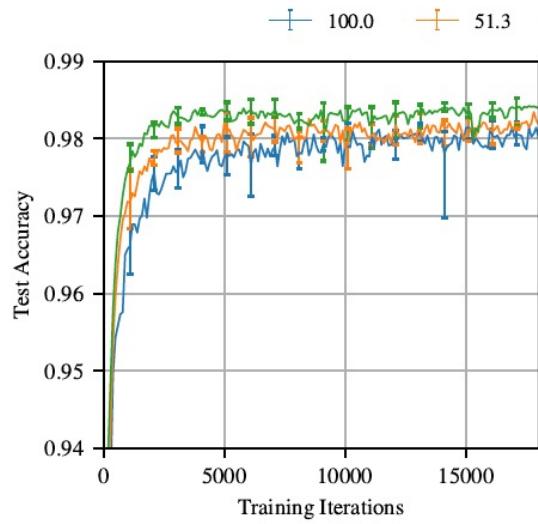
But if start from subnet work → low accuracy

bc a good initialization isn't enough to capture connectivity pattern



Empirical Evidence

Test accuracy: $21.1\% > 51.3\% \approx 7.0\% \approx 3.6\% \approx 1.9\% > 51.3 \text{ reinit} > 21.1 \text{ reinit}$
 $> 1.9\%$





Further reading

- Tishby et al. DNNs and Information Bottleneck
<https://arxiv.org/abs/1503.02406>
- Frankle & Carbin 2019 **THE LOTTERY TICKET HYPOTHESIS: FINDING SPARSE, TRAINABLE NEURAL NETWORKS**
<https://arxiv.org/pdf/1803.03635.pdf>
- DasGupta & Gupta 2002
<https://cseweb.ucsd.edu/~dasgupta/papers/jl.pdf>
- Visualizing the PHATE of neural networks
<https://arxiv.org/abs/1908.02831>