

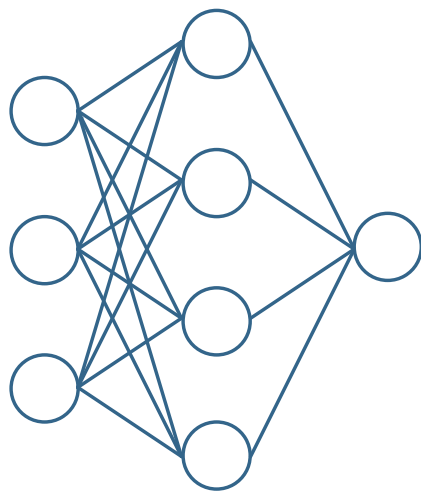
PROVING THE LOTTERY TICKET HYPOTHESIS

ICML 2020, ERAN MALACH, GILAD YEHUDAI, SHAI SHALEV-SHWARTZ, OHAD SHAMIR

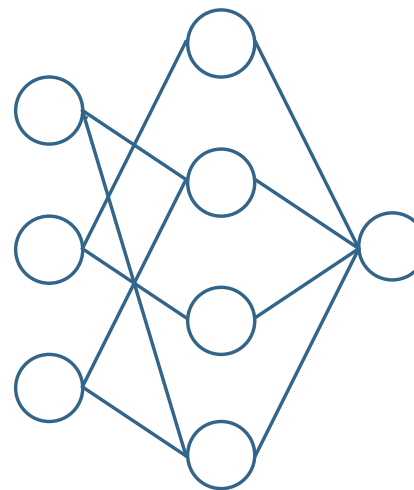
KMLR Reading Group, Geelon So, Aug. 12, 2020

PRUNING NEURAL NETS

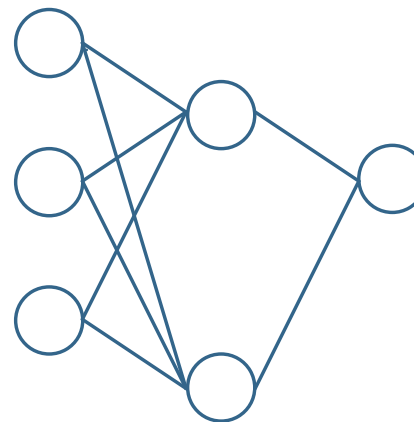
-goal: reduce space to store, time to train



weight subnetwork

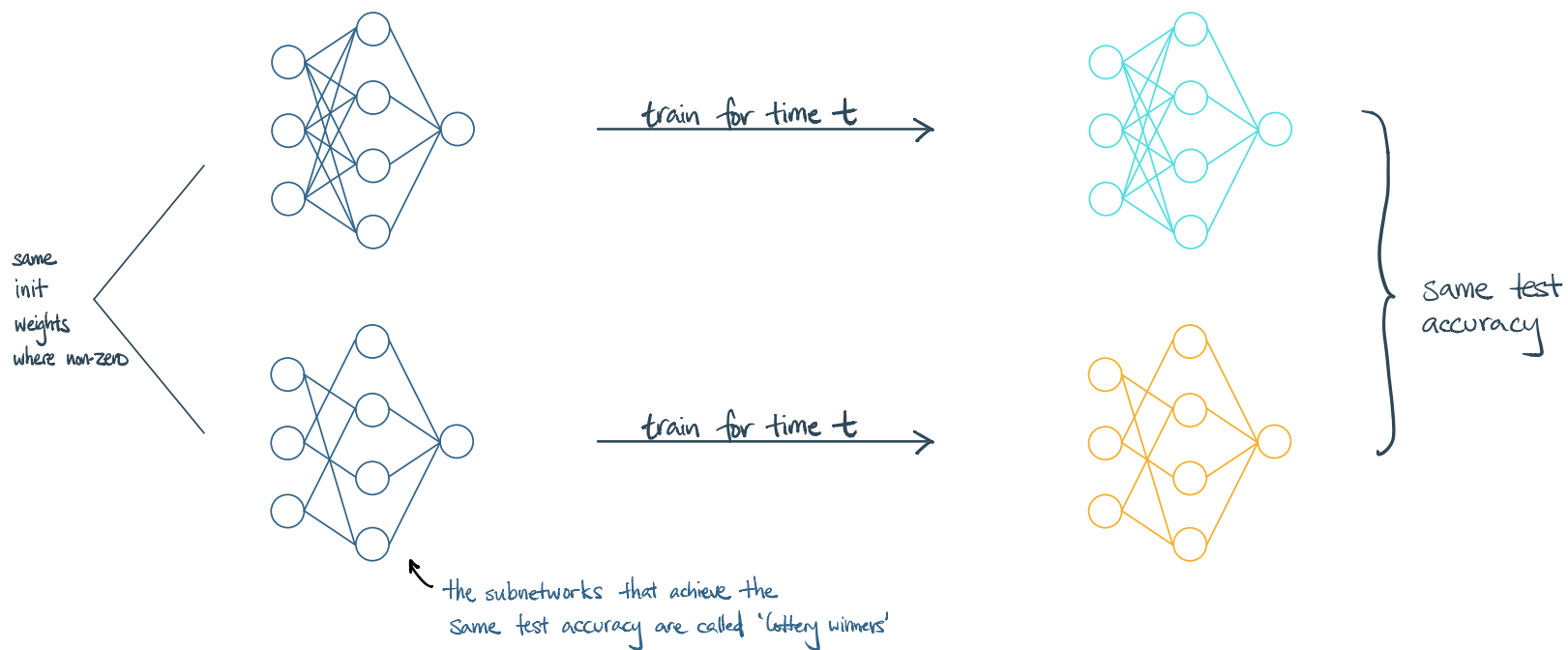


neuron subnetwork



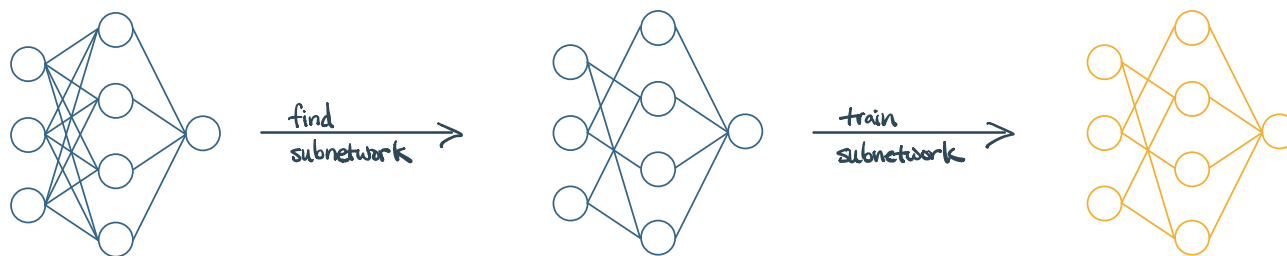
THE LOTTERY TICKET HYPOTHESIS

CONJECTURE (Frankle & Carbin 2018). A randomly-initialized, dense neural net contains a subnetwork that is initialized such that when trained in isolation, it can match the test accuracy of the original network after training for at most the same number of iterations.



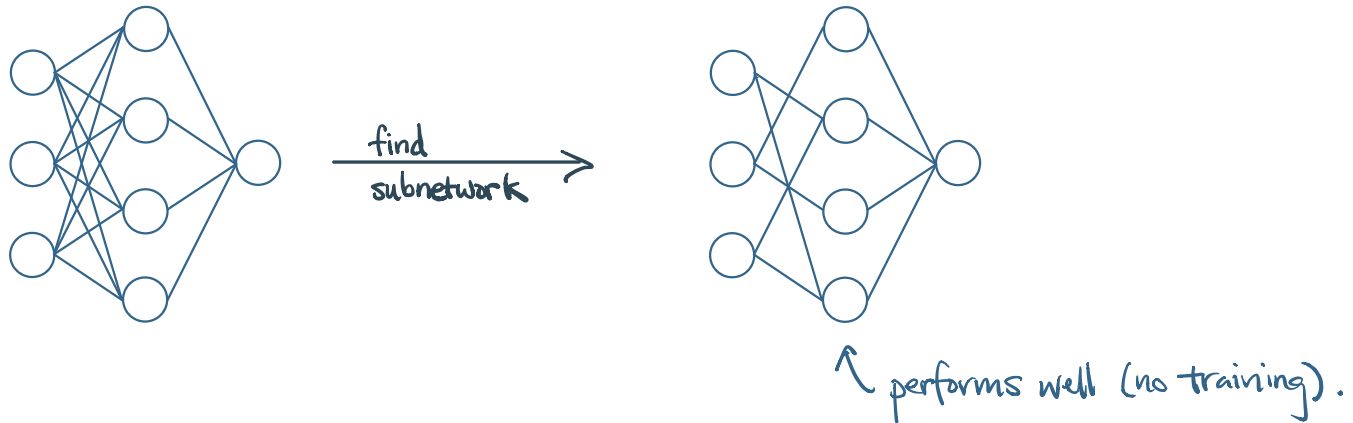
THE LOTTERY TICKET HYPOTHESIS

IMPLICATION. Instead of training a dense neural net, find suitable subnetwork to train.
i.e. it might be possible to speed up training large networks.



STRONG LOTTERY TICKET HYPOTHESIS

CONJECTURE (Ramanujan, et. al.). Within a sufficiently overparametrized neural net with random weights, there exists a subnetwork that achieves competitive accuracy.



• Proved in this paper.

RESULTS of this paper

THEOREM (weight-subnetwork, informal).

GIVENS:

· $f: \mathbb{R}^d \rightarrow \mathbb{R}$ an arbitrary ReLU neural net, depth l , width n

· $g: \mathbb{R}^d \rightarrow \mathbb{R}$ a random ReLU neural net, depth $2l$, width $\text{poly}(d, n, l, \frac{1}{\epsilon}, \log \frac{1}{\delta})$

If f satisfies mild conditions, then $\exists \tilde{g}$ a weight-subnetwork of g s.t. w.p. $1-\delta$,

$$\sup_{x \in B(0,1)} |f(x) - \tilde{g}(x)| < \epsilon.$$

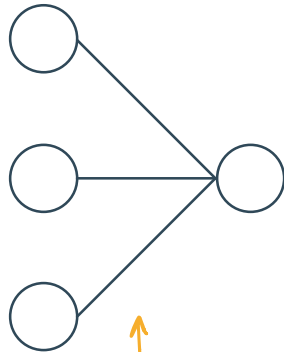
Furthermore, the number of active weights in \tilde{g} is $O(dn + n^2l)$.

↑ same order of weights in f .

PROOF IDEA

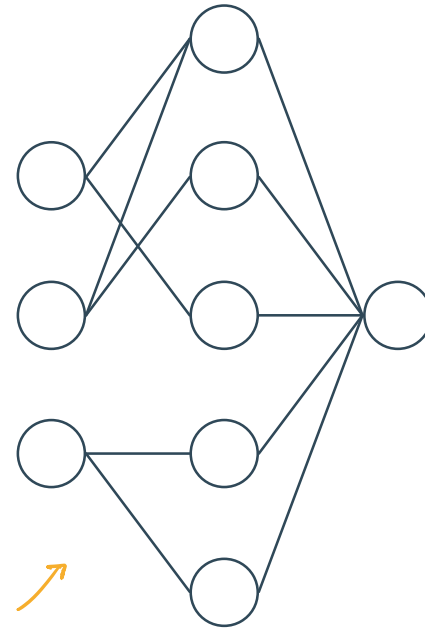
- Show that a single ReLU neuron $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is closely approximated by a weight-subnetwork of a random 2-layer net with high probability

$$f: x \mapsto \sigma\left(\sum_{i=1}^d w_i x_i\right)$$



ReLU-NN

\approx



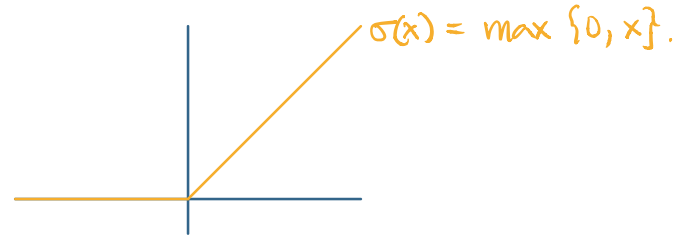
weight-subnetwork
of a random ReLU NN

PROOF IDEA

Let $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ be the ReLU activation

$$\cdot x = \sigma(x) - \sigma(-x)$$

$$\Rightarrow \sigma\left(\sum_{i=1}^d w_i x_i\right) = \sigma\left(\sum_{i=1}^d \sigma(w_i x_i) - \sum_{i=1}^d \sigma(-w_i x_i)\right)$$



Therefore, a single ReLU neuron can be written as a 2-layer ReLU NN.

PROOF IDEA

Let g be a random 2-layer neural net,

$$g(x) = \sigma \left(\sum_{j=1}^k A_j \cdot \sigma \left(\sum_{i=1}^d W_{j,i} x_i \right) \right)$$

where $A_j, W_{j,i}$ are random variables, $A \in \mathbb{R}^k$, $W \in \mathbb{R}^{k \times d}$

• Let $B \in \{0,1\}^{k \times d}$ be a bit-mask

- $B \circ W$ is a pruning of W

↑ componentwise multiplication (Hadamard product)

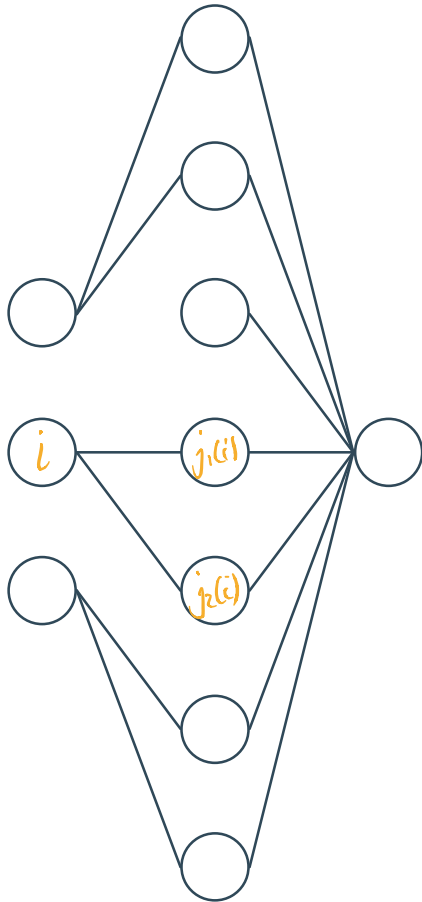
$$\bullet (B \circ W)_{ji} = \begin{cases} W_{ji} & \text{if } B_{ji} = 1 \\ 0 & \text{if } B_{ji} = 0 \end{cases}$$

PROOF IDEA

Let B be a bit-mask so that $\forall i \in [d]$,

$$j = j_1(i) \text{ or } j = j_2(i) \Leftrightarrow B_{ji} = 1$$

$$\text{and } B_{ji} = 1 \Rightarrow B_{j'i} = 0 \text{ if } j' \neq j.$$



i.e. each neuron in the first layer is connected to two unique neurons in the hidden layer.

PROOF IDEA

If B as above,

$$\begin{aligned}\tilde{g}(x) &= \sigma\left(\sum_{i=1}^d A_{j_1(i)} \cdot \sigma(W_{j_1(i),i} x_i) + \sum_{i=1}^d A_{j_2(i)} \cdot \sigma(W_{j_2(i),i} x_i)\right) \\ &= \sigma\left(\sum_{i=1}^d \text{sign}(A_{j_1(i)}) \cdot \sigma(|A_{j_1(i)}| \cdot W_{j_1(i),i} x_i) + \sum_{i=1}^d \text{sign}(A_{j_2(i)}) \cdot \sigma(|A_{j_2(i)}| \cdot W_{j_2(i),i} x_i)\right)\end{aligned}$$

Recall: $f(x) = \sigma\left(\sum_{i=1}^d \sigma(w_i x_i) - \sum_{i=1}^d \sigma(-w_i x_i)\right)$

• want $\text{sign}(A_{j_1(i)}) \neq \text{sign}(A_{j_2(i)})$

and $| |A_{j_1(i)}| \cdot W_{j_1(i),i} - \text{sign}(A_{j_1(i)}) \cdot w_i | < \frac{\varepsilon}{2d}$ (similarly for $j_2(i)$).

$\Rightarrow \sup_{x \in B(0,1)} |f(x) - \tilde{g}(x)| < \varepsilon.$

PROOF IDEA

The mild condition on $f(x) = \sigma\left(\sum_{i=1}^d w_i x_i\right)$, $w_i \in \left[-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}\right]$

w_i 's are bounded

Random initialization of $A_{j_1(i)}, w_{j_1(i),i} \sim \text{Unif}([-1,1])$

$$\Rightarrow \text{sign}(A_{j_1(i)}) \neq \text{sign}(A_{j_2(i)})$$

$$\text{and } \left| |A_{j_1(i)}| \cdot w_{j_1(i),i} - \text{sign}(A_{j_1(i)}) \cdot w_i \right| < \frac{\epsilon}{2d} \quad (\text{similarly for } j_2(i)).$$

with probability $\Omega(\epsilon/d)$.

\leadsto setting $k \gg \frac{d}{\epsilon} \Rightarrow$ for each i , we can find $j_1(i)$ and $j_2(i)$ satisfying above conditions w.h.p.

FORMAL STATEMENT

Theorem 1. Fix $\varepsilon, \delta \in (0, 1)$. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be an arbitrary ReLU neural network of depth ℓ and width n ,

$$f(x) = \sigma(w^{(\ell)} \cdot \sigma(w^{(\ell-1)} \cdot \dots \cdot \sigma(w^{(1)} \cdot x))),$$

where $w^{(1)} \in \mathbb{R}^{n \times d}$, $w^{(\ell)} \in \mathbb{R}^{1 \times n}$ and $w^{(i)} \in \mathbb{R}^{n \times n}$ for $1 < i < \ell$. Assume bounded weights, $\|w^{(i)}\|_2 \leq 1$, $\|w^{(i)}\| \leq \frac{1}{\sqrt{n_i}}$ where $n_1 = d$ and $n_i = n$ otherwise.

Let $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a random ReLU neural network with depth 2ℓ and:

$$\text{width} = \text{poly} \left(d, n, \ell, \frac{1}{\varepsilon}, \log \frac{1}{\delta} \right),$$

with weights $W^{(j)} \sim \text{Unif}([-1, 1])$ for $1 \leq j \leq 2\ell$.

With probability at least $1 - \delta$, there is a weight-subnetwork \tilde{g} of g such that:

$$\sup_{x \in B(0,1)} |\tilde{g}(x) - f(x)| \leq \varepsilon.$$

Furthermore, the number of active (non-zero) weights in \tilde{g} is $O(dn + n^2\ell)$.

} conditions on f

} conditions on g

} there is a good pruning

Proof. Apply construction for a single neuron n times for each layer, then compose layers. \square

ADDITIONAL NOTE

- the paper also shows a result for neuron-subnetworks, with a relatively simple argument reducing to random features model.

THEOREM (neuron-subnetwork, informal).

Let D be a distribution over $\mathcal{X} \times [-1, 1]$, σ be L -Lipschitz.

random features model { Let D^* be a distribution over weights $\{w: \|w\| \leq 1\}$ s.t. if $w_1, \dots, w_n \sim D^*$ w.h.p. $\exists u_1, \dots, u_n \in [-C, C]$ s.t.
$$f(x) = \sum_{i=1}^n u_i \sigma(\langle w_i, x \rangle)$$
 achieves small loss $L_D(f) \leq \epsilon$.

Let g be a random 2-layer neural network, width = $\text{poly}(C, n, L, \frac{1}{\epsilon}, \frac{1}{\delta})$ where first layer random weights $\sim D^*$, second layer $\sim \text{Unif}([-1, 1])$.

$\Rightarrow \exists \tilde{g}$ neuron-subnetwork and constant $c > 0$ s.t. $L_D(c\tilde{g}) \leq \epsilon$.

This shows the equivalence between random features models and neuron-subnetworks.

- random features models cannot efficiently approximate a single ReLU neuron.