

# One Initialization to Rule Them All?

Stefan Wijnja, Ellis Wierstra, Thomas Hamburger

January 28, 2020

Neural networks are very cool



## But they can be very demanding

- ▶ Biggest network has 16 million neurons
  - ▶ A powerful desktop needs 12 days to train 100.000 neurons
  - ▶ Would take more than 5 years to train
- ▶ Le et al (2011)
  - ▶ Cluster of 1,000 servers
  - ▶ 16,000 CPU cores
  - ▶ Ran for 3 days

## **The Lottery Ticket Hypothesis**

A randomly initialized, dense neural-network 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.

[@frankle2019]

## **The Lottery Ticket Hypothesis (redux)**

A big net contains a small net that can match the big net.

## Finding winning tickets by pruning

- ▶ Pruning
- ▶ Magnitude-based pruning
- ▶ Iterative magnitude-based pruning

# Status quo on pruning

## TODO

- ▶ [ @frankle2019 ]
  - ▶ “Winning tickets”
- ▶ *80-90% size reduction*
  - ▶ Learn faster
  - ▶ Higher test accuracy
- ▶ [ @morcoss2019 ]
  - ▶ Finding winning tickets is computationally expensive
    - ▶ Solution: possible to reuse WTs across datasets
    - ▶ Generalise without performance loss

## Our question

How do initialization algorithms compare when looking for winning tickets?

# Method

## Plan

1. Pick dataset.
2. Pick model.
3. Pick initialization methods.
4. **Write the code.**
5. **Do lots of experiments.**
6. Answer the question.

# Method

## Dataset: MNIST



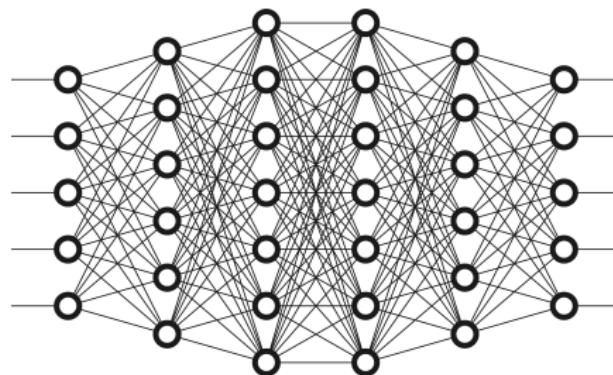
A 4x10 grid of handwritten digits. The digits are arranged in four rows. The first row contains ten '0's. The second row contains ten '1's. The third row contains ten '2's. The fourth row contains ten '3's. The fifth row contains ten '4's. The sixth row contains ten '5's. The seventh row contains ten '6's. The eighth row contains ten '7's. The ninth row contains ten '8's. The tenth row contains ten '9's. The digits are written in a cursive style and are slightly irregular.

0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9

- ▶ Handwritten digits
- ▶  $28 \times 28$  pixels
- ▶ Training set: 60k
- ▶ Test set: 10k

# Method

Model: LeNet [lecun1998]



- ▶ Fully connected
- ▶ Two hidden layers: 300 & 100 neurons → 266k weights.
- ▶ Leaky ReLU (negative slope: 0.05)

# Method

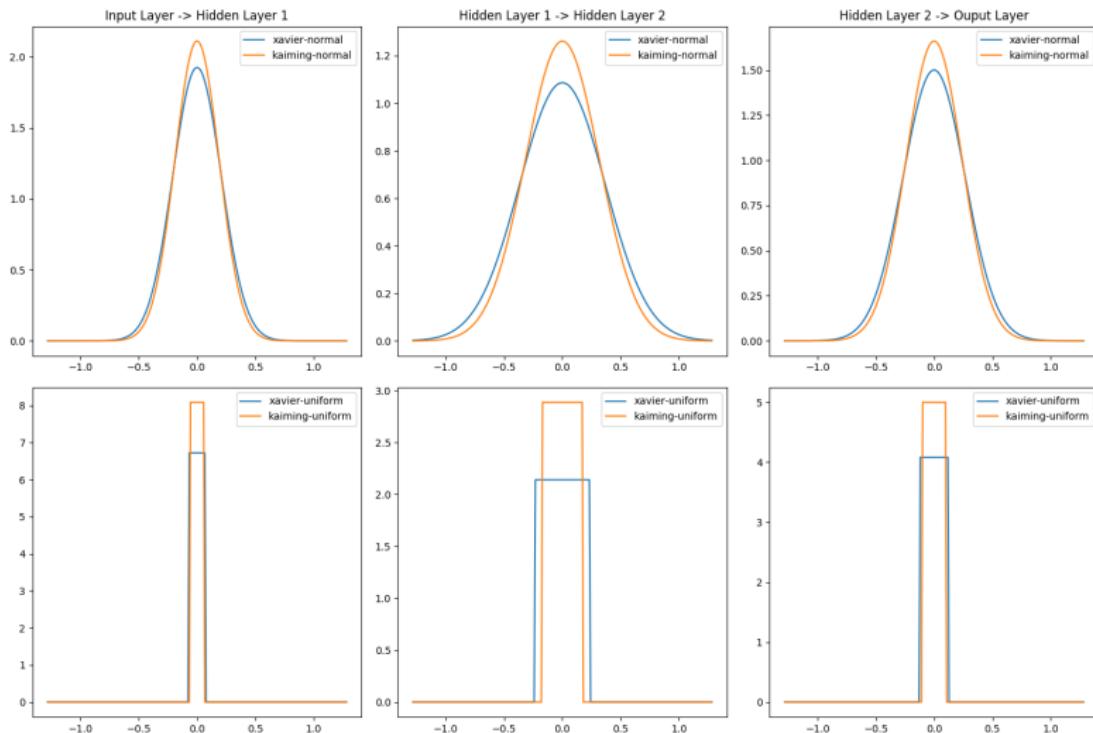
## Initialization Methods

	Xavier	Kaiming
$\mathcal{U}(-a, a)$	$\sqrt{\frac{6}{\text{fan\_in}+\text{fan\_out}}}$	$\sqrt{\frac{3}{\text{fan\_mode}}}$
$\mathcal{N}(0, \text{var})$	$\sqrt{\frac{2}{\text{fan\_in}+\text{fan\_out}}}$	$\frac{1}{\sqrt{\text{fan\_mode}}}$

- ▶  $\text{fan\_in}/\text{fan\_out}$ : number of inputs to/outputs from a neuron.  
 $\text{fan\_mode}$ :  $\text{fan\_in}$  or  $\text{fan\_out}$ .

# Method

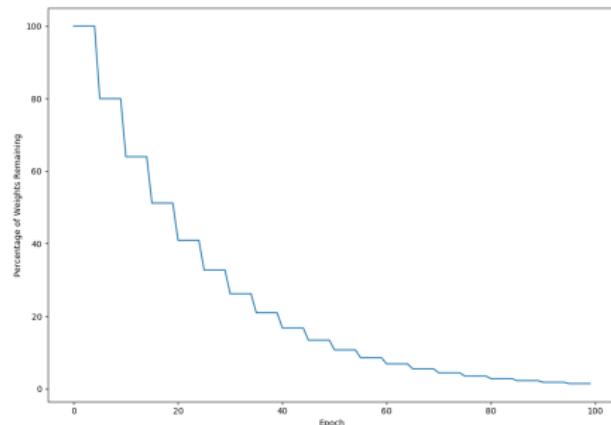
## Initialization Methods' Probability Density Functions



# Method

## Training

1. Initialize network and save weights.
2. Train for 5 epochs.
3. Trim smallest 20% of the weights.
4. Reset the weights to the saved ones.
5. GO TO 2



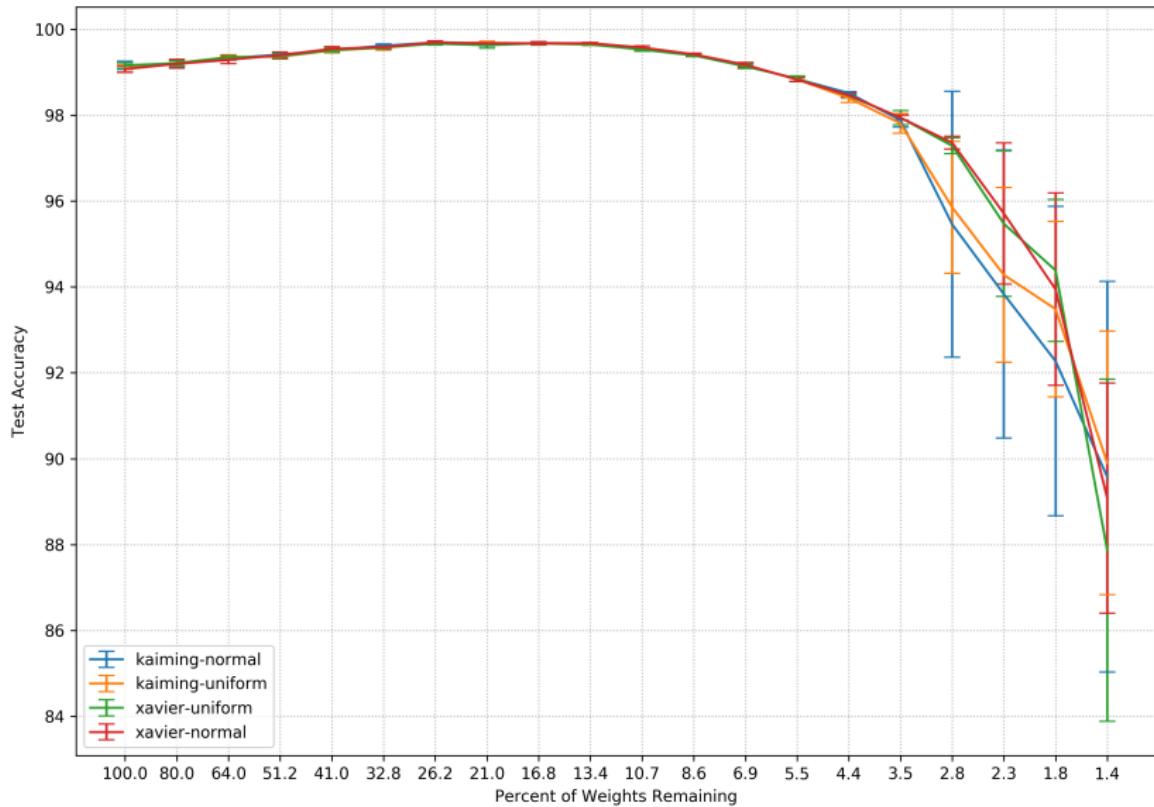
# Method

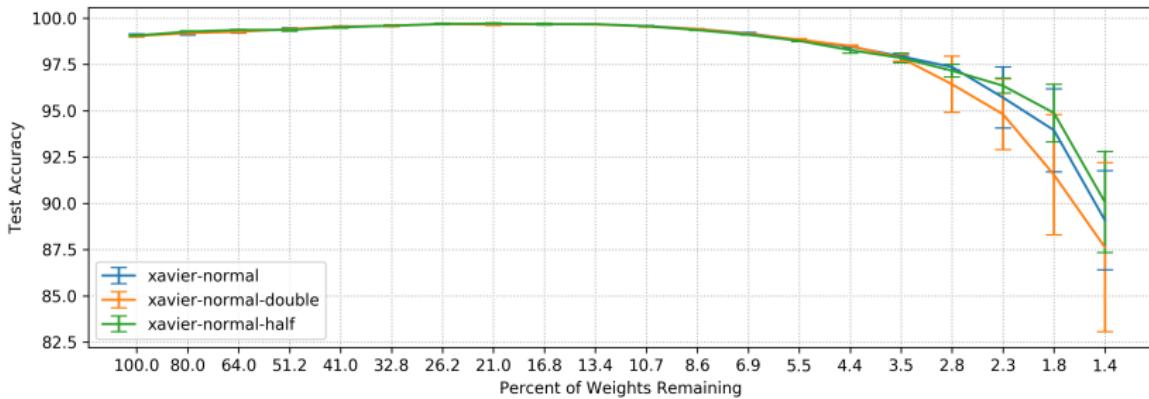
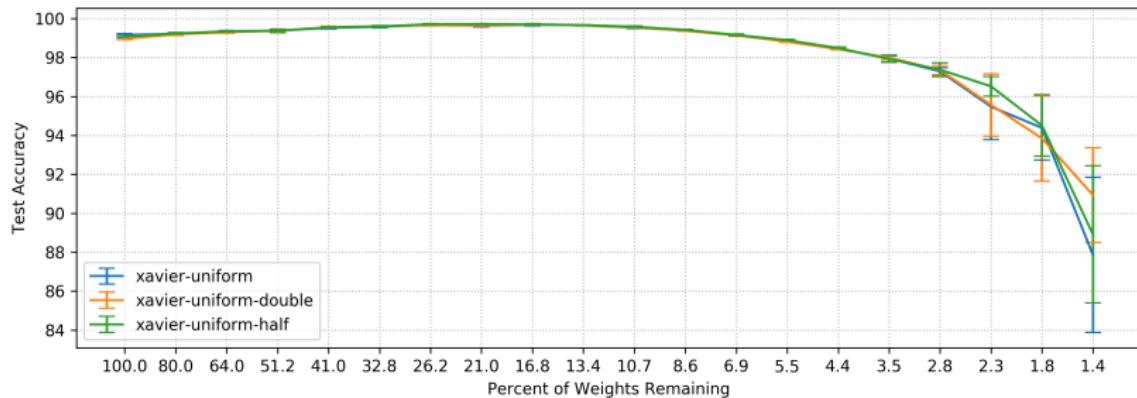
## Testing

- ▶ Save loss and accuracy on the training and testing set after every epoch.



# Results





## **Question**

How do initialization algorithms compare when looking for winning tickets?

## **Answer\***

Xavier > Kaiming

Uniform  $\approx$  Normal

Narrow > wide

## Further Research

- ▶ Grid search other datasets/models/pruning settings.
- ▶ Are pruning method and initialization method connected?

# Questions?

- ▶ Thanks
  - ▶ Andrei Apostol
  - ▶ Putri van der Linden and Lieuwe Rekker
  - ▶ BrainCreators
- ▶ Stats
  - ▶ ~599 lines of code
  - ▶ 122 training runs
- ▶ Frameworks
  - ▶ PyTorch
- ▶ Can I reproduce your results?
  - ▶ Please do. Email s.wijnja@me.com



## References