

Sparse Double Descent: Where Network Pruning Aggravates Overfitting

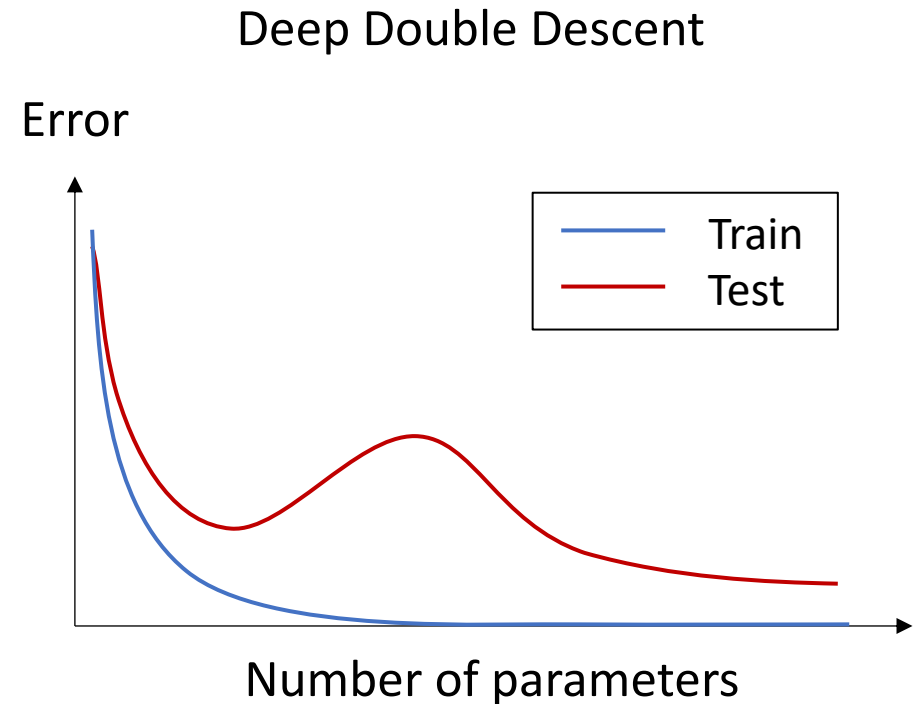
Zheng He¹, Zeke Xie^{2,3}, Quanzhi Zhu¹, Zengchang Qin¹

¹Beihang University, ²The University of Tokyo, ³RIKEN Center for AIP

Motivation

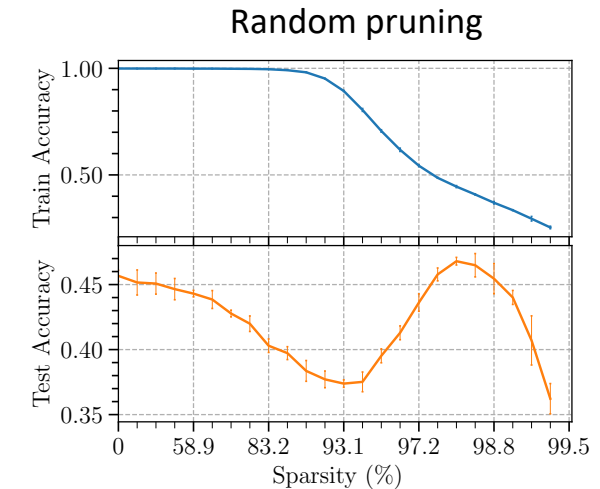
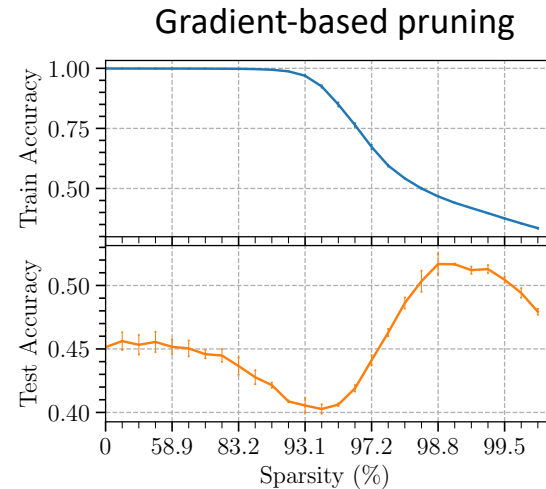
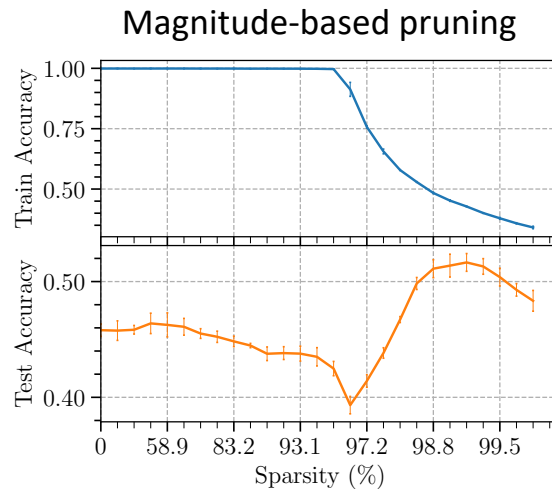
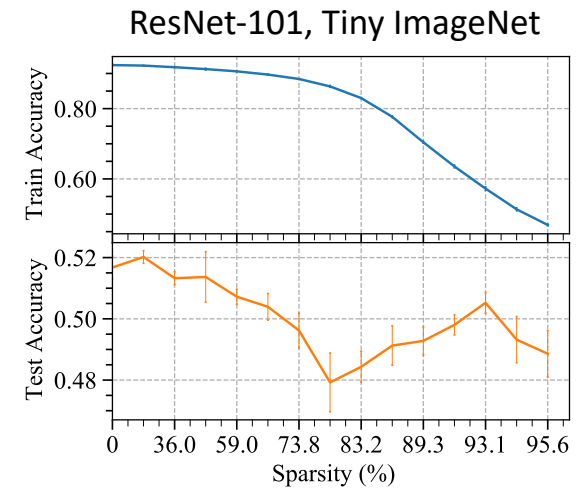
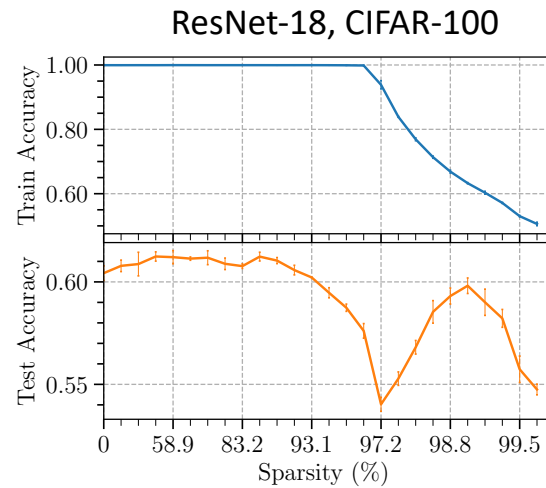
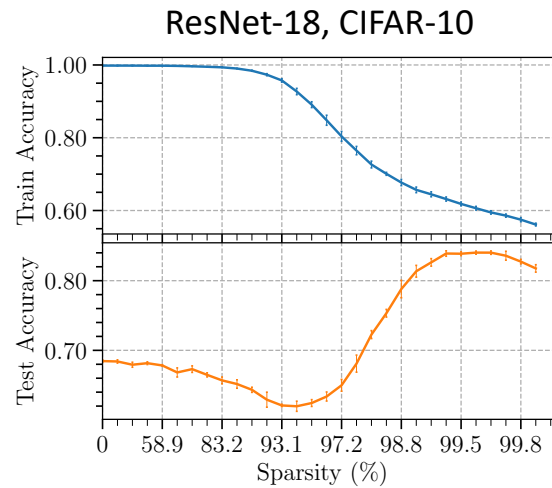
- Deep neural networks are overparameterized
- As the model capacity increases, the double descent phenomenon occurs
- Network pruning could also affect model capacity

Q: May the sparsification of DNNs also cause double descent?



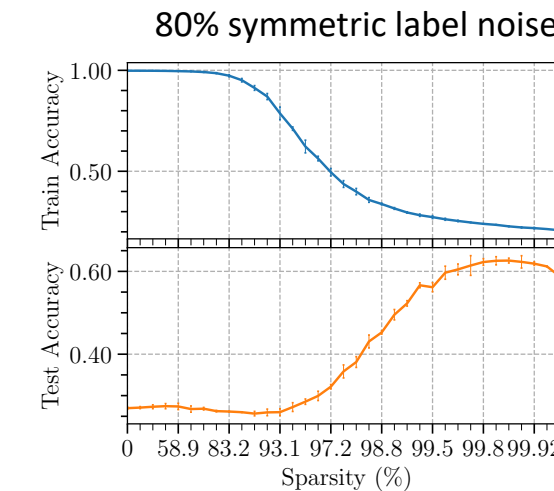
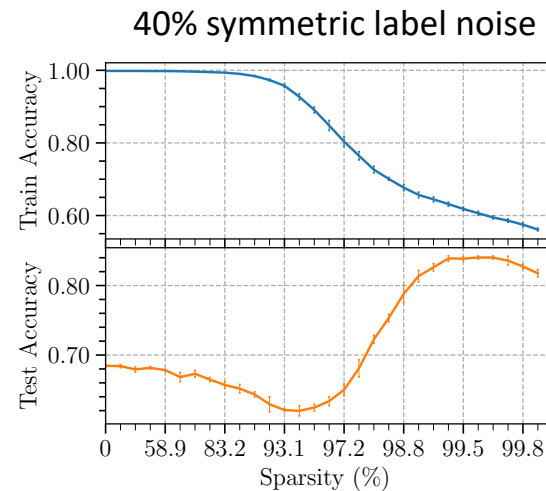
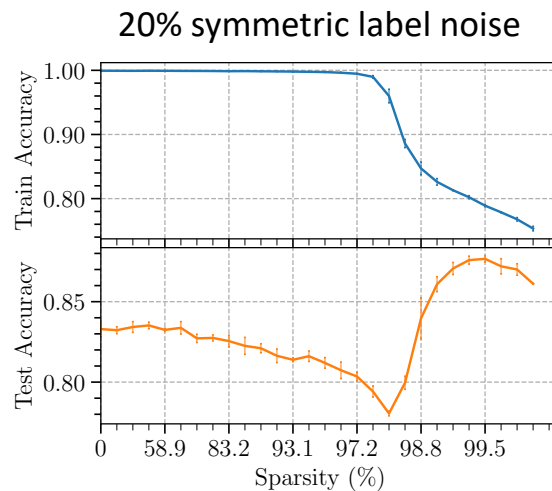
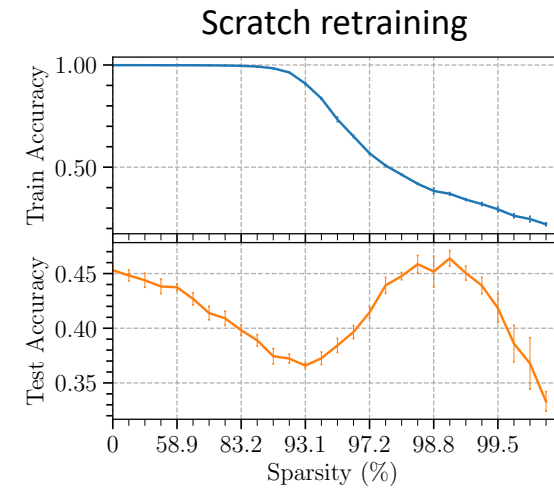
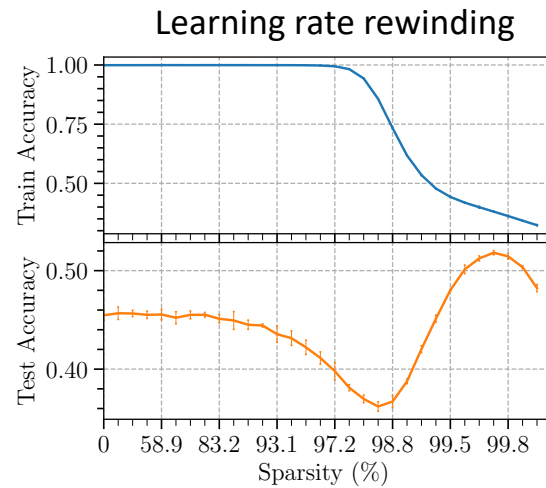
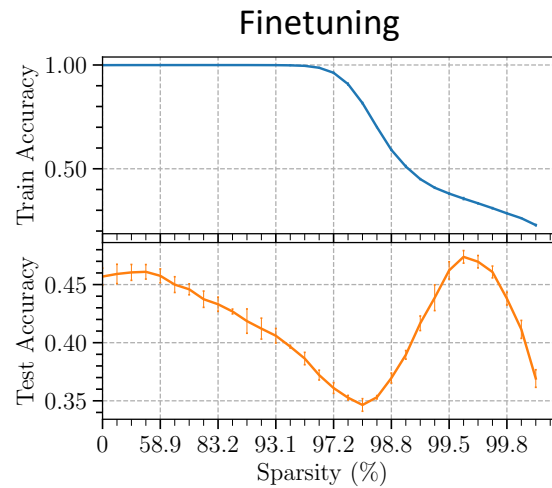
Sparse Double Descent

- Sparse double descent exists consistently across different experimental settings under label noise



Sparse Double Descent

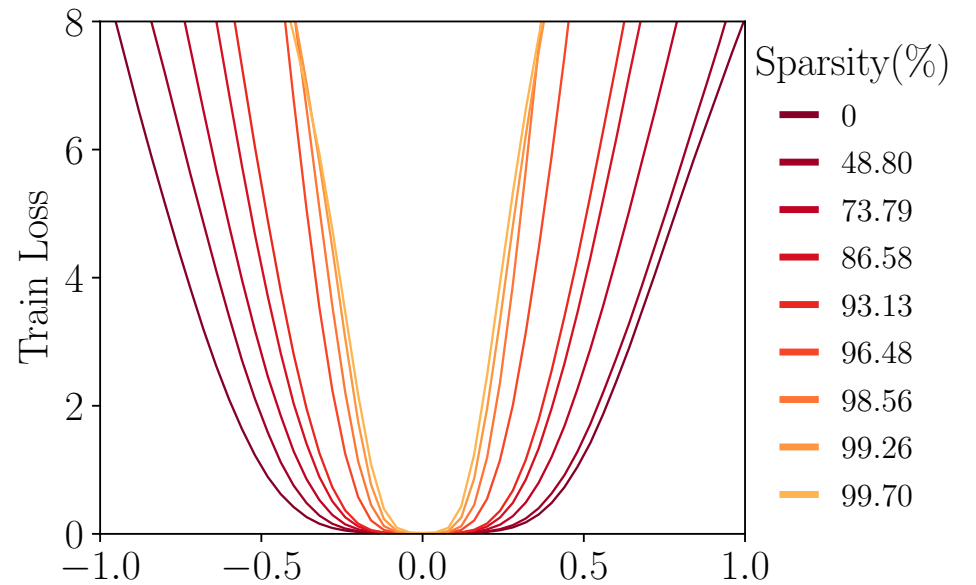
- Sparse double descent exits consistently across different experimental settings under label noise



Why Sparse Double Descent Occurs?

Minima flatness hypothesis

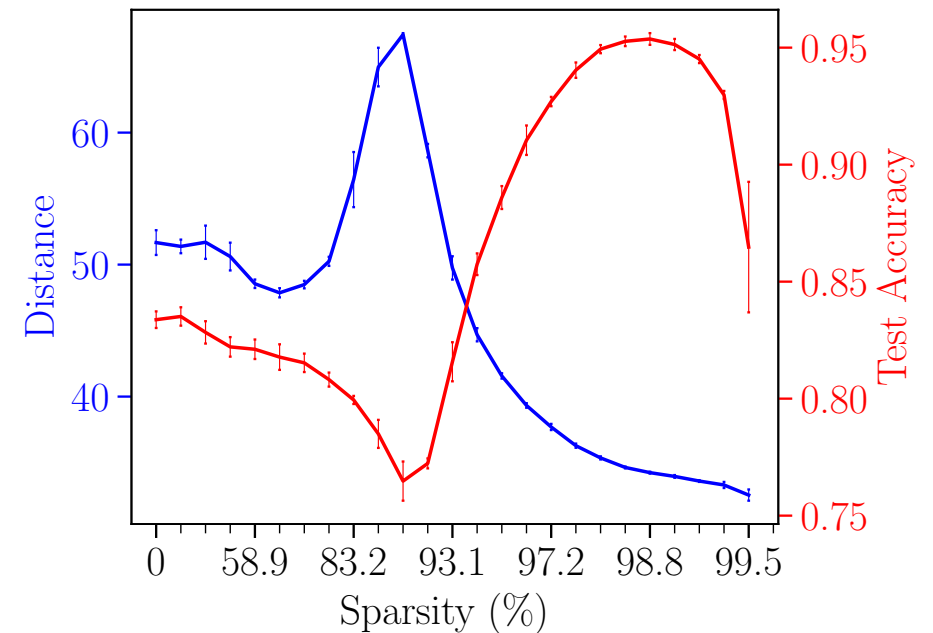
- Previous works hypothesized that pruning encourages the optimizer to move towards flatter minima [Bartoldson et al., 2020]
- Minima flatness is usually correlated with good generalization
- We observed optimizer may not converge to flat regions as sparsity increases



Why Sparse Double Descent Occurs?

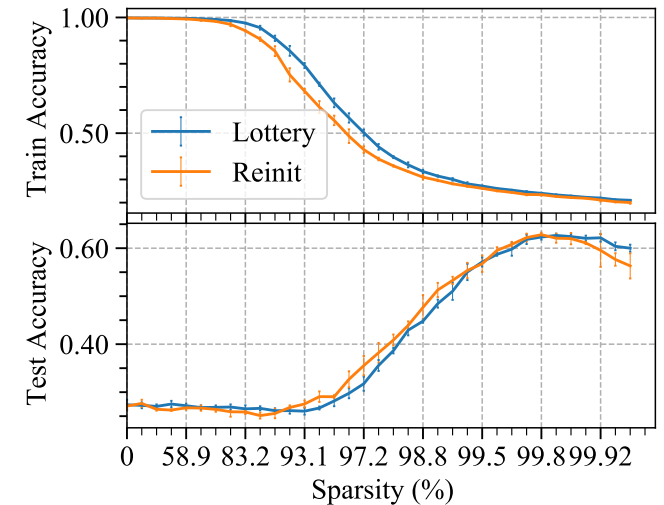
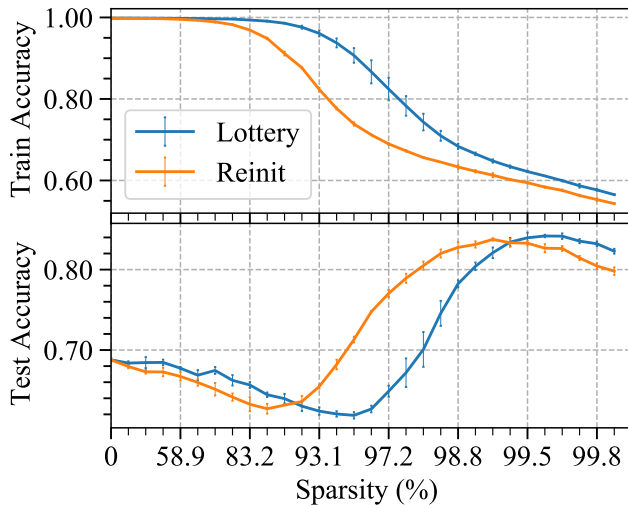
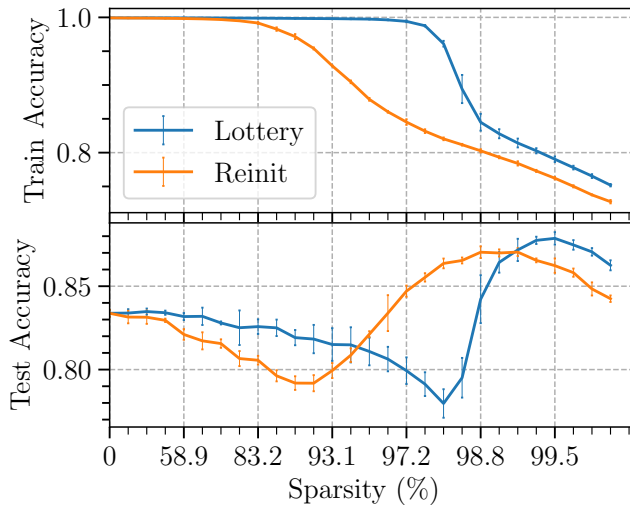
Learning Distance hypothesis

- model capacity could be restricted by the ℓ_2 learning distance from initialization [Nagarajan & Kolter, 2019]
- Pruning may affect the learning distance
- We observed the curve of learning distance correlates with test accuracy



Winning tickets may not always win

- Random reinitializations sometimes surpass the winning ticket initializations in the Lottery Ticket Hypothesis [Frankle & Carbin, 2019]



Thanks for your attention!