Bad Global Minima Exist and SGD Can Reach Them

Shengchao Liu **Dimitris Papailiopoulos Dimitris Achlioptas**

University of Athens

Rethinking Generalization Our Results Bad Global Minima Exist [Zhang et al. ICLR'17] performance on test. Overparameterized, SGD-traind models: • Can fit even completely random labels (i.e., huge SGD Can Reach Them capacity) ■ true labels • Explicit Regularization Affects Search Dynamics 2.0 random labels average loss * * shuffled pixels global minima. random pixels ← gaussian escape bad initializers. 0.5 **A Toy Example** 0.0 20 15 25 thousand steps • Yet, generalize well 1 1.0 0.8 accuracy 0.6 0.4 test(w/ aug, wd, dropout) train(w/ aug, wd, dropout) test(w/o aug, dropout) True labels Random labels True labels, train(w/o aug, dropout) 0.2 Random Init Adversarial Init Random Init test(w/o aug, wd, dropout) train(w/o aug, wd, dropout) 0.0 6000 8000 10000 4000 2000 **Confusion Data Generation** thousand training steps

Possible Explanations

- Every model that fits the training data generalizes well (No bad global minima)
- SGD "avoids" bad global minima?

Quebec Artificial Intelligence Institute (Mila), Université de Montréal

University of Wisconsin-Madison

- BAD Global Minimum: Permit fit on training data, but poor
- Very easy to construct initial conditions using only unlabeled data such that SGD converges to bad global minimum.
- Regularization helps beyond telling apart good from bad
- A combination of I2 and data augmentation allow SGD to









True labels labels, Adversarial init

Algorithm 1 Creating the Confusing (Random) Data Set **Input:** Original training dataset S; Replication factor R; Noise factor N $C = \emptyset$ for every image $x \in S$ do for i from 1 to R do $x_i \leftarrow$ zero-out a random subset comprising N% of the pixels in x $y_i \leftarrow \text{Uniformly random label}$ Add (x_i, y_i) to CTrain the architecture to 100% accuracy on C from a random initialization using vanilla SGD **Output:** The weight vector of the architecture when training ends



de Montréal







Findings and Conclusions

- Adversarial initialization causes SGD up to 40% drop in the test accuracy.
- The model found is close to the adversarial initialization.
- Data augmentation (DA), momentum (M), and I2 regularization all contribute to SGD escaping adversarial initialization.
- And two of {DA, M, I2} are enough.

Codes are available at https://github.com/chao1224/BadGlobalMinima Email: liusheng@mila.quebec