

## 1 Overview

- **Title:** Exploration and Exploitation of Deep Neural Network Loss surfaces
- **Principal Investigator:** Joan Bruna Estrach. Courant Institute, NYU, Computer Science Department, Center for Data Science and Mathematics Department. Contact: bruna@cims.nyu.edu; cell: 917-767-0617.
- **Google Sponsors:** Jascha Sohl-Dickstein, Koray Kavukcuoglu.
- **Google Contacts:** Oriol Vinyals, Vincent Vanhoucke, Jonathan Tomson.

## 2 Proposal Body

**Abstract:** The loss surface of deep neural networks has recently attracted interest in the optimization and machine learning communities as a prime example of a hard, high-dimensional, non-convex problem. Recent work has explored models from statistical physics [?], such as spin glasses [?], in order to understand the macroscopic properties of the system, but at the expense of strongly simplifying the nonlinear nature of the model.

In this proposal, we do not make any such assumption and study conditions on the data distribution and model architecture that prevent the existence of bad local minima. Together with recent results that rigorously establish that gradient descent does not get stuck on saddle points, our first objective yields guarantees that gradient descent converges to a global optimum in deep rectified networks.

The conditioning of gradient descent is the next challenge we shall address. We study this question by estimating the geometry of level sets, and we introduce an algorithm to estimate the regularity of such sets on large-scale networks. Our empirical results suggest that these sets become exponentially more curvy as the energy level decays, in accordance to what is observed in practice.

**Problem Statement:** Optimization is a critical component in Deep Learning, governing its success in different areas of computer vision, speech processing and natural language processing. The prevalent optimization strategy is Stochastic Gradient Descent, invented by Robbins and Munro in the 50s. On the one hand, the generalization of SGD is better than one could expect in generic, non-convex loss surfaces. On the other hand, one may wonder if this generic optimization strategy can be adapted to operate in the class of loss surfaces defined by deep neural networks. Some reasons to believe that such adaptation is possible come from the variety of small modifications of SGD algorithms yielding significant speedups [?, ?, ?]. This raises a number of theoretical questions as to why neural network optimization does not suffer in practice from poor local minima, as well as opportunities to leverage training data into adapted optimization.

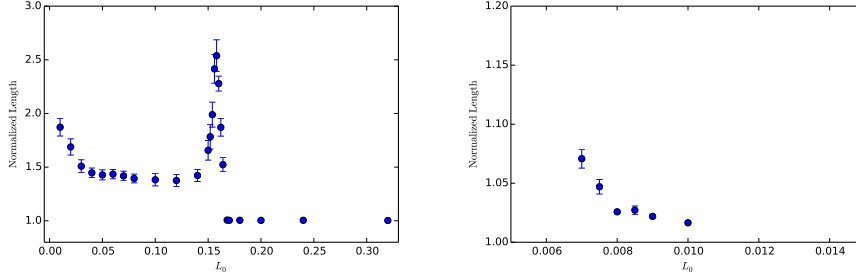


Figure 1: Normalized geodesic length as a function of the energy level for two different models.

The loss surface  $F(\theta)$  of a given model can be expressed in terms of its level sets  $\Omega_\lambda$ , which contain all parameters yielding a loss smaller or equal than the energy level  $\lambda$ :  $\Omega_\lambda = \{\theta ; F(\theta) \leq \lambda\}$ . A first question we address is about the topology of these level sets, i.e. under which conditions they are connected. This implies that one can always find a descent direction at each energy level, and therefore that no local minima exist. One can verify that linear neural networks (cascaded linear operators without pointwise nonlinearities), despite defining a non-convex loss  $F(\theta)$ , have connected level sets and therefore can be optimized globally using gradient descent strategies. The non-linear case is not true in general [?], and our first objective is to characterize the interplay between data distribution and model architecture that is needed to obtain connected level sets.

Beyond the question of whether the loss contains poor local minima, the immediate follow-up question that determines the convergence of algorithms in practice is the local conditioning of the loss surface. It is thus related not to the topology but to the shape or geometry of the level sets. As the energy level decays, one expects the level sets to reveal more complex irregular structures, which correspond to regions with small curvature.

In order to construct meaningful theoretical models, we have developed an algorithm that estimates these topological and geometrical properties on large deep architectures and datasets. Inspired by [?], given two parameter values  $\theta_1$  and  $\theta_2$  lying in the border of the same level set  $\Omega_\lambda$ , our algorithm constructs a geodesic  $\gamma(t)$  relying  $\theta_1$  to  $\theta_2$  using a dynamic programming procedure. An example of the resulting characteristics is illustrated in Figure 2.

Besides improving our understanding of stochastic optimization in the class of neural networks, our second main objective is to leverage that understanding by developing optimization strategies adapted to the class of loss surfaces at hand, i.e. learning how to learn [?]. In particular, we shall construct targets for optimization based on the geodesic paths described above: given two parameter values  $\theta_a, \theta_b$  with  $F(\theta_b) \leq F(\theta_a)$  (for example obtained during the course of a vanilla optimization method at iterations  $t$  and  $t+T$  respectively), a natural target to update  $\theta_a$  is to follow the geodesic  $\gamma(t)$  relying  $\theta_a$  to  $\theta_b$ .

The extent by which this geodesic is smooth will determine how easy it is to predict from past and present information available at  $\theta_a$ . We will explore predictive models based on recurrent neural networks, inspired from [?].

### 3 Data Policy

This proposal requires mostly unlabeled data in the form of images, speech and audio signals, and text corpus. The initial plan is to use publicly available datasets. The research outcomes will be both conference and journal papers, and also publicly available code.

### 4 Budget

- Student salary: \$50000 per year.
- Student conference/workshop travel: \$6000 per year (three conferences a year).