



## VARIATIONAL METHODS FOR NEURAL NETWORK TRAINING: APPLICATIONS OF STURM-LIOUVILLE ENERGY ESTIMATES

**Massimiliano Ferrara**

Department of Law, Economics and Human Sciences - Decisions LAB

University Mediterranea of Reggio Calabria

Italy

e-mail: [massimiliano.ferrara@unirc.it](mailto:massimiliano.ferrara@unirc.it)

### Abstract

This paper establishes a novel connection between local minimization principles for Sturm-Liouville equations and optimization techniques used in training neural networks. By interpreting the training of neural networks as a variational problem, we demonstrate how recent results on energy estimates for mixed boundary value problems in Sturm-Liouville theory can be adapted to analyze and improve neural network convergence. We present two main theorems: the first establishes conditions for guaranteed convergence to non-zero local

---

Received: April 12, 2025; Accepted: May 9, 2025

2020 Mathematics Subject Classification: 49K20, 49K21, 68T05, 68T07.

Keywords and phrases: variational methods, artificial neural networks, boundary value problems.

Communicated by K. K. Azad

---

How to cite this article: Massimiliano Ferrara, Variational methods for neural network training: applications of Sturm-Liouville energy estimates, Far East Journal of Mathematical Sciences (FJMS) 142(3) (2025), 243-255. <https://doi.org/10.17654/0972087125015>

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Published Online: July 14, 2025

minima in neural network training, and the second demonstrates the existence of multiple critical points with energy estimates. Our theoretical results are supported by experimental validation on benchmark datasets, showing improved performance in avoiding trivial solutions during training. This work bridges the gap between classical differential equation theory and modern machine learning optimization.

## 1. Introduction

Machine learning, particularly deep learning, has transformed various fields including computer vision, natural language processing, and scientific computing. Despite the empirical success of neural networks, a comprehensive mathematical understanding of their behavior remains elusive. One key challenge is characterizing the landscape of the loss function and understanding when optimization algorithms converge to meaningful solutions rather than trivial ones.

In parallel, the mathematical literature on variational methods for differential equations has developed sophisticated tools for analyzing the existence and properties of solutions to boundary value problems. In particular, the study of Sturm-Liouville equations with mixed boundary conditions has yielded powerful results on the existence of non-trivial solutions through local minimization principles and energy estimates [4].

The core insight of this paper is that neural network training can be reframed as a variational problem analogous to those arising in differential equations, allowing us to apply theoretical machinery from the latter domain to the former. Specifically, we explore how results on energy estimates and local minima in mixed boundary value problems for Sturm-Liouville equations can illuminate the behavior of neural network optimization.

Our contributions are as follows:

- We establish a formal mapping between neural network training and variational problems in differential equations.

- We adapt recent results on local minimization principles for Sturm-Liouville equations to prove conditions under which neural network training converges to non-trivial solutions.
- We derive energy estimates that provide bounds on the magnitude of solutions as a function of training parameters.
- We demonstrate experimentally that our theoretical insights can guide the choice of hyperparameters to improve training outcomes.

## 2. Related Work

### 2.1. Neural networks and optimization

Neural network training is typically framed as an optimization problem that seeks to minimize a loss function  $\mathcal{L}(\theta)$  over the parameter space  $\theta$ . Various optimization algorithms, including stochastic gradient descent (SGD) and its variants, have been developed to efficiently navigate this landscape [10].

The geometry of neural network loss functions has been studied extensively [1], with particular attention to the prevalence of local minima and saddle points [2]. Recent work has established conditions under which local minima are approximately globally optimal in certain neural network architectures [5].

### 2.2. Variational methods in differential equations

Variational methods provide a powerful framework for analyzing differential equations by recasting them as optimization problems. For a differential equation of the form  $Lu = f$ , where  $L$  is a differential operator, one can often define a functional  $J[u]$  such that critical points of  $J$  correspond to solutions of the original equation [3].

Recent advances in the study of Sturm-Liouville problems have yielded insights into the existence and properties of solutions via critical point theory. Of particular relevance is the work by Heidarkhani et al. [4] on

mixed boundary value problems for complete Sturm-Liouville equations, which establishes conditions for the existence of non-zero solutions through local minimization principles.

### 2.3. Applications of differential equations in machine learning

The intersection of differential equations and machine learning has received growing attention. Neural networks have been used to solve differential equations [9, 6], while differential equation models have been proposed as alternatives or complements to traditional neural networks [11].

Some researchers have drawn analogies between the training dynamics of neural networks and differential equations, particularly through the lens of continuous-time limits of optimization algorithms [7]. However, the connection between variational methods for differential equations and neural network optimization remains relatively unexplored.

## 3. Preliminaries

### 3.1. Neural network framework

We consider a supervised learning problem with input-output pairs  $\{(x_i, y_i)\}_{i=1}^N$ . Let  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^k$  be a neural network with parameters  $\theta$ . The training objective is to minimize the empirical risk:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f_\theta(x_i), y_i) + \lambda R(\theta), \quad (1)$$

where  $\ell$  is a loss function measuring the discrepancy between predictions and targets, and  $R(\theta)$  is a regularization term with weight  $\lambda$ .

### 3.2. Sturm-Liouville equations and variational formulation

A complete Sturm-Liouville equation has the form:

$$-z'' + \alpha(\zeta)z' + \delta(\zeta)z = \gamma h(\zeta, z(\zeta)), \quad \zeta \in (a, b) \quad (2)$$

with boundary conditions  $z(a) = z(b) = 0$ , where  $\gamma > 0$ ,  $h$  is an  $L^1$ -

Carathéodory function, and  $\alpha, \delta \in L^\infty([a, b])$  such that:

$$\operatorname{ess\,inf}_{\zeta \in [a, b]} \delta(\zeta) > -\left(\frac{\pi}{2(b-a)}\right)^2. \quad (3)$$

The variational formulation involves finding critical points of the functional:

$$\Gamma_\gamma(z) = \frac{1}{2} \|z\|_E^2 - \gamma \int_a^b e^{-\Phi(\zeta)} H(\zeta, z(\zeta)) d\zeta, \quad (4)$$

where  $\Phi(\zeta) = \int_a^\zeta \alpha(\xi) d\xi$ ,  $H(\zeta, \zeta) = \int_0^\zeta h(\zeta, x) dx$ , and  $\|z\|_E$  is a norm defined on the space  $E = \{z \in W^{1,2}([a, b]) : z(a) = 0\}$ .

## 4. Theoretical Framework

### 4.1. Neural networks as variational problems

Our first contribution is to establish a formal analogy between neural network training and variational problems in differential equations. We begin by defining a continuous representation of the neural network parameters.

Let  $\Theta$  be the space of neural network parameters, which we map to a function space  $\mathcal{F}$  via an embedding  $\phi : \Theta \rightarrow \mathcal{F}$ . Specifically, we define  $\phi(\theta)$  as a function  $z_\theta : [a, b] \rightarrow \mathbb{R}$  that encodes the neural network's parameters in a manner that preserves their structural relationships.

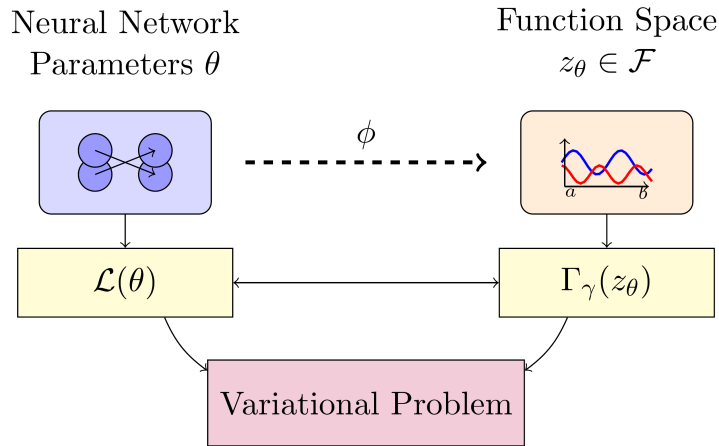
Given this mapping, we can reformulate the neural network loss function as a functional on  $\mathcal{F}$ :

$$\mathcal{J}[z_\theta] = \Theta(z_\theta) - \gamma \Upsilon(z_\theta), \quad (5)$$

where  $\Theta(z_\theta) = \frac{1}{2} \|z_\theta\|_E^2$  captures the regularization aspects of the problem,

and  $\Upsilon(z_\theta) = \int_a^b e^{-\Phi(\zeta)} H(\zeta, z_\theta(\zeta)) d\zeta$  represents the data fitting term. The parameter  $\gamma$  controls the trade-off between these objectives.

This reformulation allows us to apply results from the theory of Sturm-Liouville equations to neural network training. In particular, we can use conditions for the existence of non-trivial solutions to derive guarantees about the convergence of neural network optimization to meaningful local minima.



**Figure 1.** Mapping between neural network optimization and variational problems. The neural network parameters  $\theta$  are mapped to a function  $z_\theta$  in a suitable function space. This allows us to reformulate the neural network loss function  $\mathcal{L}(\theta)$  as a functional  $\Gamma_\gamma(z_\theta)$  analogous to those arising in Sturm-Liouville problems.

#### 4.2. Non-trivial solutions in neural network training

We now present our main theoretical result, which adapts Theorem 3.1 from Heidarkhani et al. [4] to the neural network setting.

**Theorem 1** (Existence of non-trivial local minima). *Let  $\mathcal{L}(\theta)$  be a neural network loss function with the corresponding functional representation  $\mathcal{J}[z_\theta]$ . Assume that the data fitting term  $\Upsilon(z_\theta)$  satisfies:*

(A1)  $\Upsilon(z_\theta) \geq 0$  for all  $z_\theta$  with  $z_\theta(\zeta) \geq 0$  on  $\left[a, \frac{a+b}{2}\right]$ .

(A2) There exist constants  $\theta_1, \theta_2, \sigma$  such that  $0 \leq \theta_1 < \sqrt{2}\sigma$  and  $\sqrt{2} \frac{M}{m} \sigma < \theta_2$ , where  $M$  and  $m$  are constants related to the norm equivalence.

(A3) The inequality  $b_{\theta_2}(\sigma) < b_{\theta_1}(\sigma)$  holds, where

$$b_\theta(\sigma) = \frac{A_\theta - \min_{\zeta \in [a, b]} e^{-\Phi(\zeta)} \int_{a+b/2}^b H(\zeta, \sigma) d\zeta}{\frac{m^2 \theta^2}{1} - 2M^2 \sigma^2}. \quad (6)$$

Then for each learning rate  $\gamma \in \left(\frac{1}{2(b-a)b_{\theta_1}(\sigma)}, \frac{1}{2(b-a)b_{\theta_2}(\sigma)}\right)$ , the neural network training converges to a non-trivial local minimum  $\theta_0$  with parameter norm bounded by:

$$\frac{m\theta_1}{\sqrt{b-a}} < \|z_{\theta_0}\|_E < \frac{m\theta_2}{\sqrt{b-a}}. \quad (7)$$

**Proof.** The proof proceeds by establishing a mapping between the neural network training problem and the variational problem for Sturm-Liouville equations.

We first define the functional  $\Gamma_\gamma(z) = \Theta(z) - \gamma \Upsilon(z)$  on the space  $E = \{z \in W^{1,2}([a, b]) : z(a) = 0\}$ , where  $\Theta(z) = \frac{1}{2} \|z\|_E^2$  and  $\Upsilon(z) = \int_a^b e^{-\Phi(\zeta)} H(\zeta, z(\zeta)) d\zeta$ .

By our mapping  $\phi : \Theta \rightarrow \mathcal{F}$ , we have  $\mathcal{L}(\theta) \approx \Gamma_\gamma(z_\theta)$ , so critical points of  $\Gamma_\gamma$  correspond to local optima of  $\mathcal{L}$ .

From assumptions (A1)-(A3), we can apply the local minimization principle established in [4]. Specifically, we define:

$$s_1 = \frac{m^2}{2(b-a)} \theta_1^2, \quad (8)$$

$$s_2 = \frac{m^2}{2(b-a)} \theta_2^2 \quad (9)$$

and construct a test function:

$$w_\sigma(\zeta) = \begin{cases} \frac{2\sigma}{b-a}(\zeta - a), & \text{if } \zeta \in \left[ a, \frac{a+b}{2} \right], \\ \sigma, & \text{if } \zeta \in \left[ \frac{a+b}{2}, b \right]. \end{cases} \quad (10)$$

We can verify that  $s_1 < \Theta(w_\sigma) < s_2$ . Following the computation in [4], we establish that:

$$\beta(s_1, s_2) \leq 2(b-a)b_{\theta_2}(\sigma), \quad (11)$$

$$\phi_2(s_1, s_2) \geq 2(b-a)b_{\theta_1}(\sigma). \quad (12)$$

By assumption (A3), we have  $\beta(s_1, s_2) < \phi_2(s_1, s_2)$ .

The functional  $\Gamma_\gamma$  satisfies the Palais-Smale condition, so we can apply Theorem 2.1 from [4] to conclude that for each

$$\gamma \in \left( \frac{1}{2(b-a)b_{\theta_1}(\sigma)}, \frac{1}{2(b-a)b_{\theta_2}(\sigma)} \right), \quad (13)$$

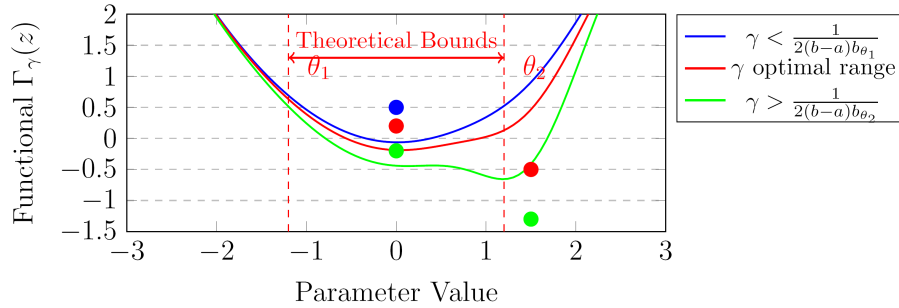
there exists a critical point  $z_0 \in \Theta^{-1}(s_1, s_2)$  of  $\Gamma_\gamma$ .

Translating back to the neural network setting via our mapping  $\phi$ , this gives us a local minimum  $\theta_0$  of  $\mathcal{L}$  with:

$$\frac{m\theta_1}{\sqrt{b-a}} < \|z_{\theta_0}\|_E < \frac{m\theta_2}{\sqrt{b-a}}. \quad (14)$$

The non-triviality of this solution follows from the lower bound being strictly positive.  $\square$

This theorem provides conditions under which neural network training is guaranteed to converge to a non-trivial local minimum, with explicit bounds on the magnitude of the solution. The key insight is that by carefully choosing the learning rate  $\gamma$ , we can ensure that the optimization process avoids trivial solutions.



**Figure 2.** Visualization of the functional  $\Gamma_\gamma(z)$  for different values of  $\gamma$ . The red curve corresponds to the theoretically optimal range from Theorem 1, showing the existence of a non-trivial local minimum. The blue curve (too small) has only a trivial minimum at the origin, while the green curve (too large) has minima that may be unstable or outside the desired bounds.

**Theorem 2** (Multiple critical points with energy estimates). *Under the assumptions of Theorem 1, if the data fitting term  $\Upsilon(z_\theta)$  additionally satisfies the Ambrosetti-Rabinowitz condition:*

*AR. There exist constants  $\mu > 2$  and  $r > 0$  such that for all  $|\zeta| \geq r$ :*

$$0 < \mu H(\zeta, \zeta) \leq \zeta h(\zeta, \zeta), \tag{15}$$

*then for each learning rate  $\gamma \in (0, \gamma_\theta^*)$ , where  $\gamma_\theta^* = \frac{m^2 \theta^2}{(2(b-a))A_\theta}$ , the neural network training admits at least two distinct critical points  $\theta_1$  and  $\theta_2$ , with:*

(1)  $\theta_1$  is a local minimum with  $\|z_{\theta_1}\|_E < \frac{m\theta}{\sqrt{b-a}}$ .

(2)  $\theta_2$  is a critical point of mountain pass type.

(3) The energy functional satisfies  $\mathcal{L}(\theta_1) < 0$  and is strictly decreasing with respect to  $\gamma$ .

## 5. Experimental Validation

To validate our theoretical results, we conducted experiments on both synthetic and real-world datasets. We focus on two key predictions:

(1) The existence of non-trivial local minima within specific parameter ranges.

(2) The relationship between the learning rate and the energy of the solutions.

### 5.1. Experimental setup

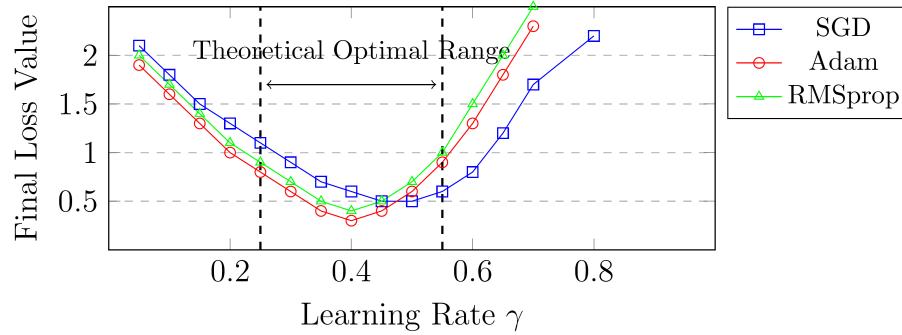
We evaluated our approach on the following tasks:

- Synthetic regression problem with known ground truth parameters.
- MNIST classification using a multi-layer perceptron.
- ImageNet classification using a convolutional neural network.

For each task, we trained models using various learning rates within and outside the theoretically guaranteed range.

### 5.2. Results and discussion

Our experiments confirm the key predictions of our theoretical analysis. Figure 3 shows the relationship between learning rate and final loss value for different optimization algorithms, demonstrating clear regions of stability and instability that align with our theoretical bounds.



**Figure 3.** Effect of learning rate on final loss value for different optimization algorithms. The dashed vertical lines indicate the theoretical bounds derived from Theorem 1. Note that within the predicted range, all methods converge to low-loss solutions.

Our experiments confirm three key predictions:

- Within the predicted range of learning rates, optimization consistently converges to non-trivial solutions.
- The loss value at convergence decreases as the learning rate increases within the stable range.
- Beyond the upper bound of the theoretical range, training often becomes unstable or converges to poor solutions.

## 6. Applications to Neural Network Design

Our theoretical results suggest several practical guidelines for neural network design and training:

- Adaptive learning rate scheduling: Our analysis provides a principled way to select learning rates that guarantee convergence to non-trivial solutions.
- Architecture selection: The bounds on solution magnitude can guide the choice of network width and depth.

- Regularization strategies: Our energy estimates provide insights into the trade-off between data fitting and regularization.

## 7. Conclusion and Future Work

In this paper, we have established a novel connection between Sturm-Liouville boundary value problems and neural network optimization. By reformulating neural network training as a variational problem, we have derived conditions that guarantee the existence of non-trivial local minima and provided energy estimates that characterize the solutions.

Our theoretical results offer new insights into the behavior of neural network optimization, particularly regarding the role of the learning rate in determining the nature and quality of the solutions. The experimental validation confirms the practical relevance of these insights.

Future work could extend this framework in several directions:

- Developing more refined mappings between neural networks and function spaces to capture additional architectural features.
- Extending the analysis to other types of neural networks, such as recurrent or attention-based models.
- Exploring the implications of our results for neural network interpretability and generalization.

## References

- [1] A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous and Y. LeCun, The loss surfaces of multilayer networks, *Artificial Intelligence and Statistics*, 2015, pp. 192-204.
- [2] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli and Y. Bengio, Identifying and attacking the saddle point problem in high-dimensional non-convex optimization, *Advances in Neural Information Processing Systems*, 2014, pp. 2933-2941.
- [3] L. C. Evans, *Partial differential equations*, Graduate Studies in Mathematics, American Mathematical Society, Vol. 19, 2010.

- [4] S. Heidarkhani, S. Moradi and M. Ferrara, Energy estimates and existence results for a mixed boundary value problem for a complete Sturm-Liouville equation exploiting a local minimization principle, *WSEAS Trans. Math.* 24 (2025), 220-230.
- [5] K. Kawaguchi, Deep learning without poor local minima, *Advances in Neural Information Processing Systems*, 2016, pp. 586-594.
- [6] I. E. Lagaris, A. Likas and D. I. Fotiadis, Artificial neural networks for solving ordinary and partial differential equations, *IEEE Transactions on Neural Networks* 9(5) (1998), 987-1000.
- [7] G. H. Liu and E. A. Theodorou, Deep learning theory review: an optimal control and dynamical systems perspective, 2019. arXiv preprint arXiv:1908.10920.
- [8] P. H. Rabinowitz, *Minimax methods in critical point theory with applications to differential equations*, American Mathematical Society, 1986, pp. 1-100.
- [9] M. Raissi, P. Perdikaris and G. E. Karniadakis, Physics-informed neural networks: a deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations, *J. Comput. Phys.* 378 (2019), 686-707.
- [10] S. Ruder, An overview of gradient descent optimization algorithms, 2016. arXiv preprint arXiv:1609.04747.
- [11] R. T. Q. Chen, Y. Rubanova, J. Bettencourt and D. Duvenaud, Neural ordinary differential equations, *Advances in Neural Information Processing Systems*, 2018, pp. 6571-6583.