

Elliptic Loss Regularization

Ali Hasan^{1,2} Haoming Yang¹ Yuting Ng¹ Vahid Tarokh¹

¹Duke University ²Morgan Stanley

Motivating Questions

1. Can we anticipate the behavior of the loss within regions of low training data density?
2. Can we enforce this using a scalable computational solution that applies to general function approximators?

Empirical Risk Minimization (ERM) does not explicitly regularize model behavior in regions of low training data density. A common way to transformed risk minimization (TRM) with \mathcal{T} being randomly sampled transformation parameterized by ϕ

$$\min_{\theta} \mathbb{E}_{\delta_{\mathcal{X} \times \mathcal{Y}}}[\ell(f_{\theta}(X), y)] \quad (\text{ERM})$$

$$\min_{\theta} \mathbb{E}_{\phi} \mathbb{E}_{\delta_{\mathcal{X} \times \mathcal{Y}}}[\ell(f_{\theta}(\mathcal{T}_X^{\phi} X), \mathcal{T}_y^{\phi} y)] \quad (\text{TRM})$$

We propose a particular class of transformations to enforce the loss landscape to solve a partial differential equation (PDE) with useful properties.

Our Contributions

1. We describe a new regularization scheme that enforces the model loss as a function of inputs to satisfy an elliptic PDE;
2. We theoretically characterize the practical properties of this regularization through PDE theory;
3. We introduce an efficient computational approach that endows the properties of the elliptic regularization.

Why Elliptic PDEs?

PDEs characterize the joint rates in change of different variables. We apply PDE to study how the loss changes as a function of perturbations in the input space, this:

- provides insights into the **robustness** of the model for regions outside of the training data;
- provides a correspondence between the model loss as a function of inputs and qualitative properties of elliptic PDEs to **bound the loss**.

We define the *loss landscape* as a function $u(X, y) : \mathcal{D} \rightarrow \mathbb{R}_+$. Our goal is to prescribe the function $u(X, y)$ with a specified level of regularity over the data space in a way that also imbues $\ell(f_{\theta}(X), y)$ with regularity and thereby obtain desirable properties.

$$\begin{aligned} \min_{\theta} u(X, y), \quad (X, y) \in \mathcal{D} \\ 0 = \sigma \nabla^2 u(X, y), \quad (X, y) \in \mathcal{D} \\ u(X, y) = \ell(f_{\theta}(X), y), \quad (X, y) \in \partial \mathcal{D} \end{aligned} \quad \begin{matrix} (1) \\ (2) \end{matrix}$$

A Scalable Solution

PDE in (1) can be solved by the following expectation with a by applying the Feynman-Kac formula, which requires sampling a stopping time τ :

$$u(x, y) = \min_{\theta} \mathbb{E}[\ell(f_{\theta}(x_{\tau}), y_{\tau}) | x_0 = x, y_0 = y] \quad (3)$$

To avoid infinite stopping times, we instead take expectations over Brownian bridges connecting points on the boundary:

$$\min_{\theta} \mathbb{E}_{(X_{\rightarrow}, y_{\rightarrow}) \times (X_{\leftarrow}, y_{\leftarrow}) \sim P(\Pi)} \mathbb{E}_{X_s, y_s \sim \text{BB}_{X_{\rightarrow}, y_{\rightarrow}}^{X_{\leftarrow}, y_{\leftarrow}}} \left[\int_0^1 \ell(f_{\theta}(X_s), y_s) ds \right] \quad (4)$$

for all s , where we denote $\text{BB}_{X_{\rightarrow}, y_{\rightarrow}}^{X_{\leftarrow}, y_{\leftarrow}}$ as a Brownian bridge sample path where \rightarrow denotes starting points and \leftarrow denotes end points and Π is the set of points in the support of $\delta_{(\mathcal{X} \times \mathcal{Y})} \times \delta_{(\mathcal{X} \times \mathcal{Y})}$.

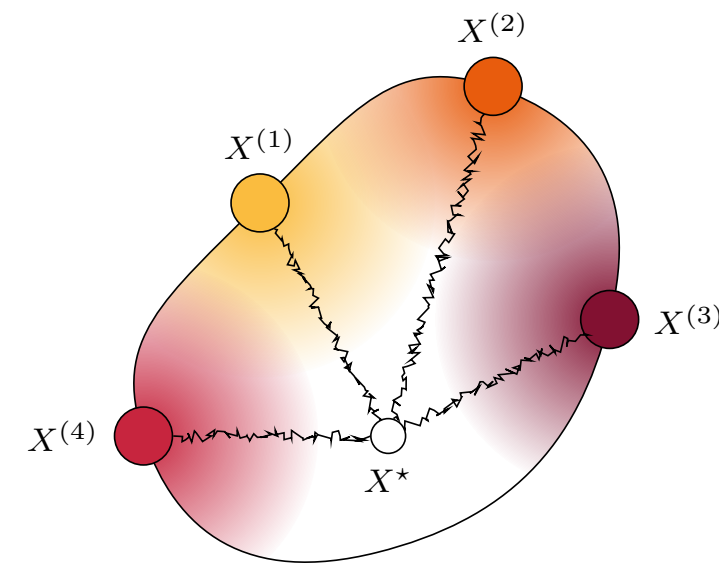


Figure 1. Illustration of the loss values over a domain with 4 points on the boundary. The expected loss at point X^* is composed of losses at ε -balls around $X^{(i)}$, $i = 1 \dots 4$. Black paths represent sample paths starting at X^* .

The Maximum Principle

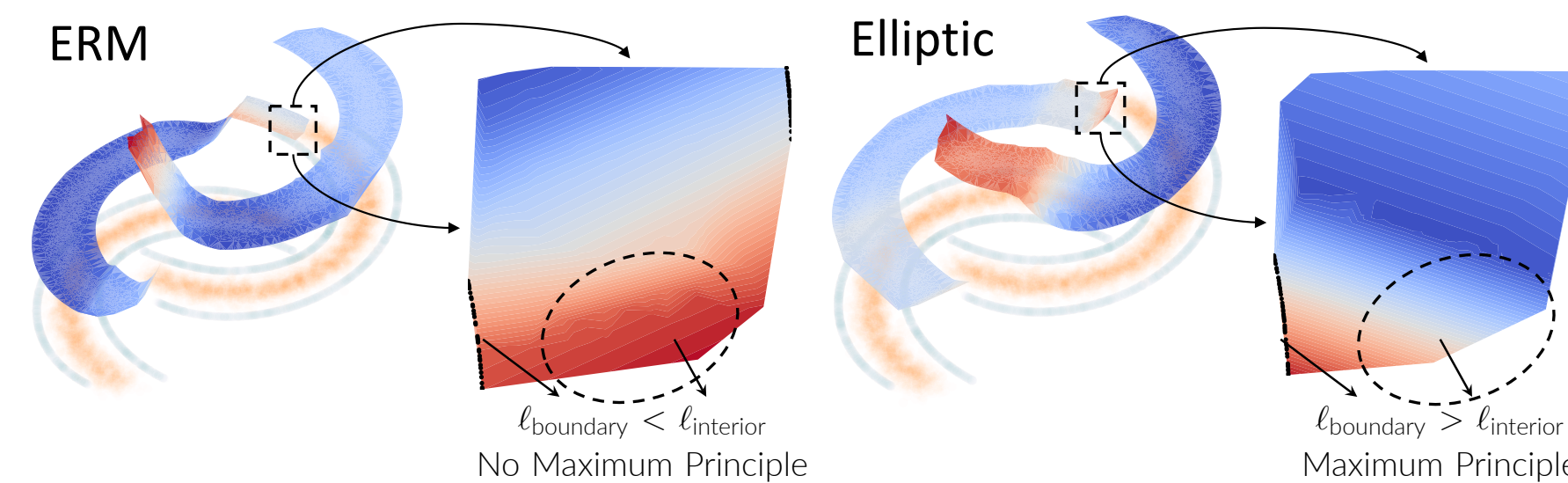
Once the optimization is solved, the loss landscape $u(x, y)$ should follow:

The Maximum Principle Consider any point $X, y \in \mathcal{D}$, suppose the function pairs u, f_{θ} solves (1). Then, the expected loss u at X, y satisfies the inequality:

$$\min_{X, y \in \mathcal{X} \times \mathcal{Y}} \ell(f_{\theta}(X), y) \leq u(X, y) \leq \max_{X, y \in \mathcal{X} \times \mathcal{Y}} \ell(f_{\theta}(X), y).$$

To illustrate this (Figure 2) consider two-moon classifiers, one trained with ERM and another with elliptic regularization, using only the blue scattered boundary data and evaluate on the red interior data to plot the loss surface.

Zooming in on the interior of the ERM loss surface, the training loss exceeds the loss of boundary (circled area) whereas elliptic regularization bounds the loss, as the theory predicts.



Empirical Evaluation

We apply elliptic regularization on classic group-imbalance, data-shift, and spurious correlation datasets and compare the results for classification:

Algorithm	WaterBirds		CelebA		Camelyon17
	Avg(%)	Worst(%)	Avg(%)	Worst(%)	
CVaR-DRO [Levy et al., 2020]	90.3 \pm 1.2	77.2 \pm 2.2	86.8 \pm 0.7	76.9 \pm 3.1	70.5 \pm 5.1
CVaR-DORO [Zhai et al., 2021]	91.5 \pm 0.7	77.0 \pm 2.8	89.6 \pm 0.4	75.6 \pm 4.2	67.3 \pm 7.2
χ^2 -DRO [Levy et al., 2020]	88.3 \pm 1.5	74.0 \pm 1.8	87.7 \pm 0.3	78.4 \pm 3.4	68.0 \pm 6.7
χ^2 -DORO [Zhai et al., 2021]	89.5 \pm 1.0	76.0 \pm 3.1	87.0 \pm 0.6	75.6 \pm 3.4	68.0 \pm 7.5
Elliptic + IW	92.0 \pm 0.3	84.1 \pm 1.1	91.3 \pm 0.3	77.4 \pm 4.5	77.9 \pm 3.0
Two-stage: JTT [Liu et al., 2021]	93.6 \pm NA	86.0 \pm NA	88.0 \pm NA	81.1 \pm NA	69.1 \pm 6.4
Two-stage: UMIX [Han et al., 2022]	93.0 \pm 0.5	90.0 \pm 1.1	90.1 \pm 0.4	85.3 \pm 4.1	75.1 \pm 5.9

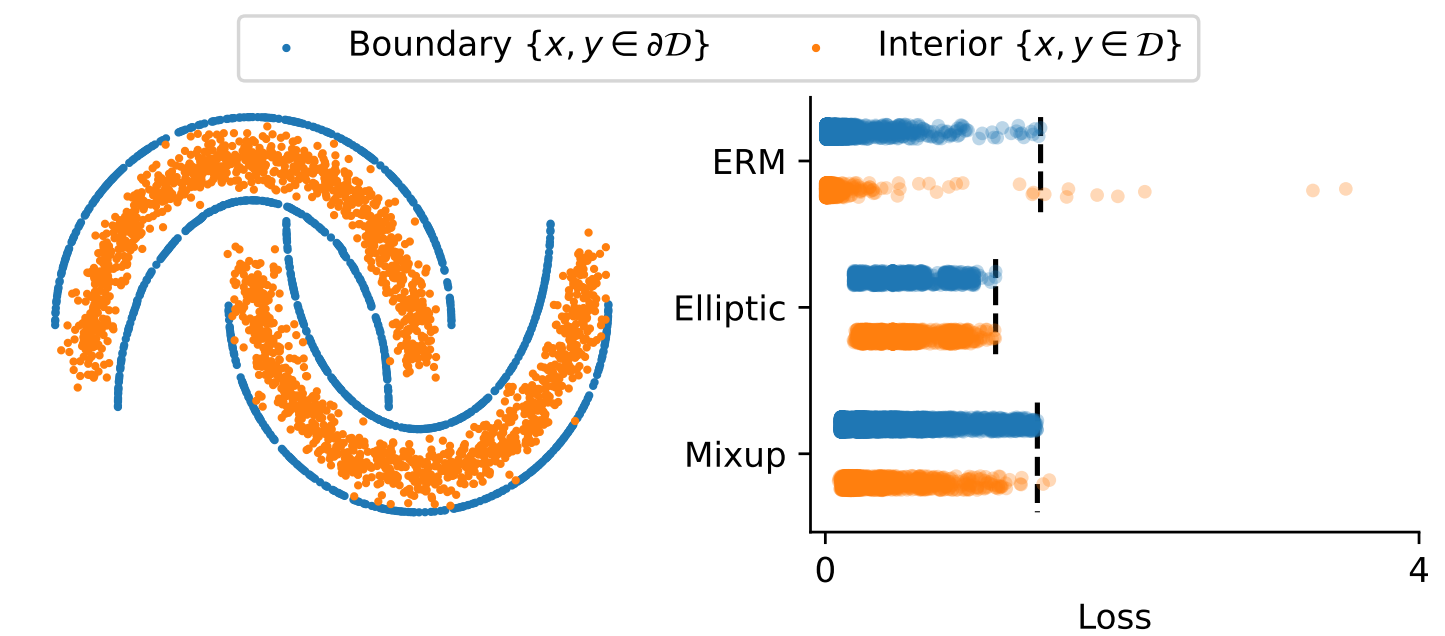
for robust regression under (sub)domain shift and spurious correlation:

Algorithm	SkillCraft		Crime		RCF-MNIST	
	Avg	Worst	Avg	Worst	Avg	Worst
C-mixup [Yao et al., 2022]	6.27 \pm 0.537	8.83 \pm 1.010	0.132 \pm 0.003	0.167 \pm 0.010	0.165 \pm 0.001	0.180 \pm 0.001
Elliptic	5.97 \pm 0.283	9.17 \pm 1.150	0.132 \pm 0.003	0.164 \pm 0.01	0.162 \pm 0.002	0.178 \pm 0.002

and with scenarios with additional 50% label corruption:

Method	Breast		Blood		Path		OrganC	
	Avg(%)	Worst(%)	Avg(%)	Worst(%)	Avg(%)	Worst(%)	Avg(%)	Worst(%)
ERM	80.0 \pm 3.4	32.4 \pm 13.8	81.4 \pm 2.7	56.3 \pm 17.3	55.7 \pm 3.8	10.7 \pm 8.6	84.1 \pm 1.9	65.1 \pm 8.2
mixup	83.1 \pm 2.4	47.1 \pm 10.9	78.8 \pm 3.1	41.7 \pm 19.3	54.3 \pm 2.1	1.2 \pm 2.1	85.6 \pm 2.4	56.8 \pm 11.0
mixupE	76.5 \pm 2.8	16.7 \pm 11.3	74.2 \pm 3.6	31.9 \pm 16.9	52.5 \pm 2.0	2.7 \pm 7.6	70.4 \pm 7.6	35.6 \pm 16.8
Elliptic	87.6 \pm 1.8	67.6 \pm 4.4	85.5 \pm 1.6	60.2 \pm 14.9	62.3 \pm 2.5	26.4 \pm 10.2	87.7 \pm 0.9	65.9 \pm 3.4
+ IW	87.3 \pm 0.9	67.6 \pm 3.0	84.1 \pm 4.0	68.0 \pm 13.7	62.7 \pm 1.8	20.9 \pm 6.8	88.0 \pm 0.7	66.2 \pm 4.8

Figure 2. Maximum principle satisfied for the elliptic regularization when trained on boundary data.



References

- Z. Han, Z. Liang, F. Yang, L. Liu, L. Li, Y. Bian, P. Zhao, B. Wu, C. Zhang, and J. Yao. Umix: Improving importance weighting for subpopulation shift via uncertainty-aware mixup. *Advances in Neural Information Processing Systems*, 35:37704–37718, 2022.
- D. Levy, Y. Carmon, J. C. Duchi, and A. Sidford. Large-scale methods for distributionally robust optimization. *Advances in Neural Information Processing Systems*, 33:8847–8860, 2020.
- E. Z. Liu, B. Haghighi, A. S. Chen, A. Raghunathan, P. W. Koh, S. Sagawa, P. Liang, and C. Finn. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, pages 6781–6792. PMLR, 2021.
- H. Yao, Y. Wang, L. Zhang, J. Y. Zou, and C. Finn. C-mixup: Improving generalization in regression. *Advances in Neural Information Processing Systems*, 35:3361–3376, 2022.
- R. Zhai, C. Dan, Z. Kolter, and P. Ravikumar. Doro: Distributional and outlier robust optimization. In *International Conference on Machine Learning*, pages 12345–12355. PMLR, 2021.