Bootstrap Consistency Regularization for Stable Neural Network Predictions^{*}

Gaurav Sood[†]

July 16, 2025

Abstract

Neural networks exhibit substantial prediction variability when retrained on bootstrap samples of the same dataset, undermining reliability in deployment scenarios requiring consistent decision boundaries. We propose a bootstrap-aware regularization technique that directly minimizes prediction variance across data resamples during training. Our method simultaneously trains multiple shadow copies of a network, each on bootstrap resamples of mini-batches, while penalizing disagreement between their predictions. Empirical evaluation on tabular datasets demonstrates 25–80% reductions in bootstrap prediction variance with accuracy degradation limited to one percentage point. Unlike existing stability approaches that target weight-space curvature or optimization noise, our method directly optimizes the quantity of practical interest: prediction consistency under data resampling.

1 Introduction

The deployment of neural networks in production systems requires not only predictive accuracy but also consistency across model updates. When a model is retrained on fresh samples from the same distribution, predictions on identical inputs should remain stable within the bounds justified by sampling uncertainty. However, neural networks commonly exhibit substantial *refit variance*—the phenomenon whereby predictions vary significantly when models are trained on different bootstrap samples of the training data.

This instability poses significant challenges across multiple domains. In production machine learning systems, model updates may reverse binary classifications on borderline cases, creating inconsistent user experiences. Scientific applications require stable models for fair method comparisons and reliable bootstrap-based confidence intervals. Regulated

^{*}https://github.com/finite-sample/consistentshade.

[†]Gaurav can be reached at gsood07@gmail.com

industries face compliance issues when prediction variability across training runs triggers audit procedures.

Existing approaches to neural network stability primarily target indirect proxies for the desired behavior. Sharpness-Aware Minimization Foret et al. (2021) and related methods penalize weight-space curvature under the assumption that flatter minima correspond to more stable predictions. Stochastic regularization techniques such as R-Drop Liang et al. (2021) control prediction consistency under network noise but do not address data resampling variance. Teacher-student methods like Mean Teacher Tarvainen and Valpola (2017) stabilize optimization dynamics while remaining agnostic to bootstrap variance.

We propose a fundamentally different approach: *bootstrap-aware regularization* that directly minimizes prediction variance across data resamples. Our method trains multiple shadow copies of a model simultaneously, each processing bootstrap resamples of training mini-batches, while explicitly penalizing disagreement between their predictions. This approach directly targets the quantity of interest rather than relying on indirect proxies.

2 Method

2.1 Problem Formulation

Consider a training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ drawn from distribution P, and a parameterized model $f_{\theta} : \mathcal{X} \to \mathcal{Y}$. For any test input $x \in \mathcal{X}$, we define the *bootstrap prediction variance* as:

$$\sigma_{\text{boot}}^2(x) = \mathbb{E}_{\mathcal{D}' \sim \text{Boot}(\mathcal{D})} \left[(\hat{f}_{\mathcal{D}'}(x) - \mathbb{E}[\hat{f}_{\mathcal{D}'}(x)])^2 \right]$$
(1)

where $\hat{f}_{\mathcal{D}'}$ denotes the model obtained by training on bootstrap sample \mathcal{D}' , and Boot(\mathcal{D}) represents the bootstrap distribution over datasets of size n sampled with replacement from \mathcal{D} .

Our objective combines standard empirical risk minimization with explicit bootstrap variance regularization:

$$\min_{\theta} \mathbb{E}_{(x,y)\sim P}[\ell(f_{\theta}(x), y)] + \lambda \cdot \mathbb{E}_{x\sim P_{X}}[\sigma_{\text{boot}}^{2}(x)]$$
(2)

where ℓ is a loss function, P_X is the marginal distribution of inputs, and $\lambda > 0$ controls the regularization strength.

2.2 Bootstrap-Aware Training Algorithm

Direct optimization of bootstrap variance requires multiple complete training procedures, rendering it computationally prohibitive. We approximate this objective using *micro-bootstrap* resampling within mini-batches, enabling efficient joint optimization.

Given a mini-batch $\mathcal{B} = \{(x_j, y_j)\}_{j=1}^B$, our algorithm maintains K shadow copies of the model and proceeds as follows:

- 1. Micro-bootstrap resampling: For each shadow model $k \in \{1, ..., K\}$, generate bootstrap indices $idx^{(k)} = \{i_1^{(k)}, \ldots, i_B^{(k)}\}$ where each $i_j^{(k)} \sim Uniform(\{1, \ldots, B\})$ independently with replacement.
- 2. Shadow predictions: Compute predictions for each shadow model on its bootstrap resample:

$$\mathbf{p}^{(k)} = f_{\theta^{(k)}}(\mathbf{x}_{\mathrm{idx}^{(k)}}) \in \mathbb{R}^B$$
(3)

3. Joint objective optimization: Minimize the combined loss:

$$\mathcal{L} = \frac{1}{K} \sum_{k=1}^{K} \frac{1}{|\mathrm{idx}^{(k)}|} \sum_{i \in \mathrm{idx}^{(k)}} \ell(f_{\theta^{(k)}}(x_i), y_i) + \lambda \cdot \frac{1}{B} \sum_{j=1}^{B} \mathrm{Var}_k[p_j^{(k)}]$$
(4)

where $\operatorname{Var}_{k}[p_{j}^{(k)}] = \frac{1}{K} \sum_{k=1}^{K} (p_{j}^{(k)} - \bar{p}_{j})^{2}$ and $\bar{p}_{j} = \frac{1}{K} \sum_{k=1}^{K} p_{j}^{(k)}$.

The algorithm updates all shadow models jointly using shared gradient information, encouraging consensus across bootstrap resamples while maintaining individual adaptation to each resample's characteristics.

2.3 Implementation Considerations

Computational overhead: Training requires K forward passes per mini-batch, increasing computational cost by a factor of approximately K. Memory requirements scale linearly with K due to the need to store multiple model copies.

Inference: At test time, predictions can be obtained from any single shadow model or their ensemble average. No additional computational cost is incurred during inference compared to standard training.

Hyperparameter selection: We fix K = 3 and $\lambda = 0.05$ across all experiments based on preliminary validation studies. These values provide a reasonable balance between stability improvement and computational overhead.

3 Experimental Setup

3.1 Datasets and Tasks

We evaluate our approach on four tabular datasets spanning regression and binary classification:

• Synthetic Regression: 20-dimensional Gaussian features with quadratic target function, n = 1000

- California Housing: Median house value prediction, 8 features, n = 20,640
- Adult Income: Binary income classification, 14 features, n = 48,842
- German Credit Risk: Binary credit risk classification, 20 features, n = 1000

All datasets employ stratified 75%/25% train/test splits to ensure representative evaluation sets.

3.2 Model Architecture and Training

We employ a standardized two-layer multilayer perceptron architecture across all experiments:

- Input layer to 64 hidden units with ReLU activation and dropout (p = 0.1)
- Hidden layer to 128 units with ReLU activation and dropout (p = 0.1)
- Output layer (1 unit for regression, 2 for classification)

Training configuration includes Adam optimization with learning rate 10^{-3} , batch size 64, and 25 epochs. We conduct 30 independent training runs per experimental condition to ensure statistical reliability.

3.3 Evaluation Metrics

Predictive performance: We report test RMSE for regression tasks and classification accuracy for binary tasks.

Bootstrap stability: For each test input x_i , we compute the sample variance of predictions across 30 independent model fits and summarize stability as:

StabilityRMSE =
$$\sqrt{\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \operatorname{Var}_{\text{fit}}[\hat{f}(x_i)]}$$
 (5)

This metric quantifies prediction variability in the same units as the target variable, facilitating interpretation across tasks.

4 Results

4.1 Main Experimental Results

Table 1 presents our primary experimental findings comparing standard empirical risk minimization against bootstrap-aware training with K = 3 shadow models and $\lambda = 0.05$.

The results demonstrate substantial improvements in bootstrap stability across all datasets, with reductions in prediction variance ranging from 26% to 81%. Importantly,

Dataset	Metric	Baseline	Bootstrap-Aware	Δ Error	Δ Stability
Synthetic	RMSE	23.88 ± 0.37	29.46 ± 0.64	+23%	-38%
California	RMSE	0.591 ± 0.005	0.598 ± 0.004	+1%	- 26%
Adult	Accuracy	0.826 ± 0.001	0.825 ± 0.001	-0.1pp	-81%
German Credit	Accuracy	0.697 ± 0.009	0.688 ± 0.006	-0.9pp	-48%

Table 1. Comparison of predictive performance and bootstrap stability. Values represent mean \pm standard deviation across 30 independent runs.

these stability gains come at minimal cost to predictive accuracy: real-world datasets (California Housing, Adult Income, German Credit) show accuracy degradation limited to one percentage point or less.

The synthetic regression task represents a challenging scenario where model capacity substantially exceeds data complexity, leading to higher baseline instability. Even in this worst-case setting, our method achieves a 38% reduction in bootstrap variance, albeit with a more substantial accuracy cost (+23% RMSE).

4.2 Stability-Accuracy Trade-off Analysis

To quantify the practical significance of our stability improvements, we analyze the decomposition of total prediction uncertainty. For the Adult Income dataset, bootstrap variance accounts for 29% of total prediction uncertainty under standard training, reducing to 10% with bootstrap-aware regularization. This represents a meaningful reduction in the uncertainty attributable to training procedure variability rather than fundamental task difficulty.

5 Related Work

Our approach differs fundamentally from existing stability methods in directly targeting prediction variance under data resampling rather than indirect proxies.

Sharpness-based methods such as Sharpness-Aware Minimization Foret et al. (2021) and Entropy-SGD Chaudhari et al. (2017) seek flatter loss surfaces under the hypothesis that such minima correspond to more stable predictions. However, the relationship between weight-space geometry and bootstrap prediction variance remains theoretically unclear.

Stochastic consistency methods like R-Drop Liang et al. (2021) enforce agreement between predictions under different dropout masks, addressing network stochasticity but not data resampling variance. These methods are architecture-specific and do not generalize to bootstrap stability.

Teacher-student approaches including Mean Teacher Tarvainen and Valpola (2017) stabilize training dynamics through exponential moving averages of model weights. While effective for reducing optimization noise, these methods do not explicitly address variability under data resampling.

Distributionally robust optimization methods such as χ^2 -DRO Duchi and Namkoong (2021) penalize loss variance across data subsets. However, loss variance does not directly correspond to prediction variance—models may achieve similar loss values while producing substantially different predictions.

6 Discussion and Limitations

6.1 Computational Considerations

The primary limitation of our approach is computational overhead. Training K shadow models increases memory requirements by a factor of K and training time by approximately $2-3\times$ due to additional forward and backward passes. For large-scale models or datasets, this overhead may prove prohibitive.

Future work could address this limitation through influence function approximations that enable single-model estimation of bootstrap variance, eliminating the need for multiple shadow models while preserving the direct optimization objective.

6.2 Theoretical Understanding

While our empirical results demonstrate clear benefits, the theoretical relationship between micro-bootstrap variance within mini-batches and full bootstrap variance across complete datasets merits further investigation. Establishing formal conditions under which our approximation remains valid would strengthen the theoretical foundations of the approach.

6.3 Hyperparameter Sensitivity

Our experiments employ fixed hyperparameters (K = 3, $\lambda = 0.05$) across all datasets. While these values prove effective in our evaluation, optimal settings may vary with model architecture, dataset characteristics, and task requirements. Developing principled approaches for hyperparameter selection represents an important direction for future research.

7 Conclusion

We have presented a bootstrap-aware regularization technique that directly addresses prediction instability under data resampling, a fundamental challenge in reliable machine learning deployment. Our method achieves substantial reductions in bootstrap prediction variance (25–80%) while maintaining competitive predictive accuracy across tabular datasets.

The key insight underlying our approach is that stability under data resampling can be effectively improved by explicitly penalizing prediction disagreement across bootstrap resamples during training, rather than relying on indirect proxies such as weight-space curvature or optimization dynamics. While computationally more demanding than single-model training, our method provides a direct solution to a pervasive problem in machine learning reliability.

Future research directions include developing computationally efficient approximations through influence functions, establishing theoretical guarantees for the micro-bootstrap approximation, and extending the approach to other neural architectures and domains beyond tabular data.

References

- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. (2017). Entropy-SGD: Biasing gradient descent into wide valleys. In *International Conference on Learning Representations*.
- Duchi, J. C. and Namkoong, H. (2021). Learning models with uniform performance via distributionally robust optimization. *Annals of Statistics*, 49(3):1378–1406.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. (2021). Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Repre*sentations.
- Liang, X., Wu, L., Li, J., Wang, Y., Meng, Q., Qin, T., Chen, W., Zhang, M., and Liu, T.-Y. (2021). R-drop: Regularized dropout for neural networks. In Advances in Neural Information Processing Systems, volume 34, pages 10890–10905.
- Tarvainen, A. and Valpola, H. (2017). Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In Advances in Neural Information Processing Systems, volume 30, pages 1195–1204.