

# Stochastic Covariance-Based Initialization (SCBI): A Scalable Warm-Start Strategy for High-Dimensional Linear Models

Fares Ashraf  
*Independent Researcher*  
*farsashraf44@gmail.com*

February 10, 2026

## Abstract

Training high-dimensional linear and logistic regression models on massive datasets typically relies on iterative optimization algorithms (e.g., Stochastic Gradient Descent) initialized with random weights. While effective, this “cold-start” approach ignores available statistical properties of the data, resulting in unnecessary computational expense during the early epochs of convergence. In this paper, we propose **Stochastic Covariance-Based Initialization (SCBI)**, a non-iterative method that approximates the optimal solution for convex linear models prior to training. By leveraging GPU parallelism, our method computes the marginal correlation of features against the target and adjusts for feature interactions via a block-wise approximation of the inverse covariance matrix. We mitigate the computational cost of  $O(d^3)$  matrix inversion for high-dimensional feature spaces ( $d > 2000$ ) through a bagging strategy. Experimental results on synthetic and real-world benchmarks (California Housing, Forest Cover Type) demonstrate that SCBI achieves near-optimal accuracy in near-constant time, reducing the required training epochs by up to 95% compared to standard He and Xavier initialization.

**Code Availability:** The official implementation is available at:  
<https://github.com/fares3010/SCBI>

## 1 Introduction

The initialization of neural network weights is a critical factor in the convergence speed and final performance of machine learning models. Standard initialization strategies, such as Xavier (Glorot) and He initialization, are designed to preserve signal variance across layers, preventing vanishing or exploding gradients [1, 2]. However, these methods are **semantically blind**: they initialize weights based on the *dimensions* of the data, not the *content* of the data.

For linear layers, logistic regression, and the classification heads of transfer learning models, the optimal weights are mathematically deterministic. However, calculating the closed-form solution (Normal Equation) involves inverting the feature interaction matrix  $(X^T X)^{-1}$ , an operation with  $O(Nd^2 + d^3)$  complexity. For high-dimensional datasets ( $d > 10^4$ ), this is computationally prohibitive, forcing practitioners to rely on slow, iterative gradient descent starting from random noise.

We propose **Stochastic Covariance-Based Initialization (SCBI)**, a hybrid approach that bridges the gap between closed-form statistics and iterative deep learning. SCBI utilizes GPU parallelism to compute a stochastic approximation of the inverse covariance matrix via bagging (bootstrap aggregating). This provides a “warm start” that places parameters in the immediate convex basin of the global minimum.

Our contributions are as follows:

1. We define the **SCBI Algorithm**, a GPU-accelerated method to approximate second-order statistical relationships for initialization.
2. We introduce a linear-complexity **Correlation Damping heuristic** as a lightweight alternative to matrix inversion.
3. We demonstrate empirically that SCBI reduces initial loss by over **60%** on high-dimensional synthetic data and effectively “solves” regression tasks before the first training epoch.

## 2 Methodology

### 2.1 Problem Formulation

Consider a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  where  $x_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}^c$ . We aim to initialize the weights  $W \in \mathbb{R}^{d \times c}$  and bias  $b \in \mathbb{R}^c$  of a linear map  $f(x) = W^T x + b$ .

Standard initialization draws  $W \sim \mathcal{N}(0, \sigma^2)$ . Our goal is to estimate  $W_{init}$  such that the loss  $\mathcal{L}(W_{init})$  approximates the global minimum  $\mathcal{L}(W^*)$ .

### 2.2 The SCBI Algorithm (Matrix Formulation)

Instead of solving the exact Normal Equation on the full dataset, we employ a bagging strategy. Let  $\mathcal{B} = \{B_1, \dots, B_K\}$  be  $K$  random subsets of the data. For each subset, we construct the augmented feature matrix  $\tilde{X}_k$  (including a bias column) and target matrix  $Y_k$ .

The estimator for the initialization weights  $\theta_{SCBI}$  is defined as the ensemble average of the Ridge-regularized solutions for each subset:

$$\theta_{SCBI} = \frac{1}{K} \sum_{k=1}^K \left[ \left( \tilde{X}_k^T \tilde{X}_k + \lambda I \right)^{-1} \tilde{X}_k^T Y_k \right] \quad (1)$$

Where:

- $\tilde{X}_k^T \tilde{X}_k \in \mathbb{R}^{(d+1) \times (d+1)}$  is the **Interaction (Covariance) Matrix**.
- $\tilde{X}_k^T Y_k$  represents the **Marginal Correlations** (Slopes).
- $\lambda$  is a regularization term (Ridge) to ensure numerical stability against multicollinearity.

By computing these terms on GPU, we exploit the hardware’s optimized GEMM (General Matrix Multiply) capabilities, rendering the calculation time negligible compared to a single training epoch.

### 2.3 The Correlation Damping Approximation

For extremely high-dimensional spaces where matrix inversion is infeasible even on subsets, we derive a diagonal approximation. We determine the weight  $w_i$  for feature  $i$  by damping its univariate slope by the sum of its absolute correlations with all other features  $j$ :

$$w_i \approx \frac{\text{Slope}(x_i, y)}{1 + \sum_{j \neq i} |\text{Corr}(x_i, x_j)|} \quad (2)$$

**Why Correlation?** We explicitly utilize the Pearson correlation coefficient rather than raw regression slope for the interaction term (denominator). This ensures the penalty factor is:

1. **Scale Invariant:** Unaffected by unit differences (e.g., meters vs. kilometers).
2. **Symmetric:**  $\text{Corr}(i, j) = \text{Corr}(j, i)$ .
3. **Bounded:** The absolute value ensures the denominator is strictly  $\geq 1$ , preventing numerical explosion.

### 3 Classical Regression Experimental Setup

We validate SCBI against standard random initialization (He/Xavier) using the Adam optimizer.

**Datasets:**

- **Synthetic High-Dim:** 5,000 samples, 2,000 features (only 50 informative).
- **California Housing:** Standard regression benchmark (8 features).
- **Forest Cover Type:** Large-scale multi-class classification (54 features, 7 classes, 581k samples).

**Implementation:** All experiments were conducted using PyTorch with CUDA acceleration. SCBI was configured with  $K = 10$  subsets and  $\lambda = 1.0$ .

## 4 Results and Discussion (Classical)

### 4.1 Convergence Analysis

Table 1 summarizes the initial loss (Epoch 0) for both methods. SCBI consistently provides a “Warm Start,” effectively skipping the initial exploration phase of optimization.

Dataset	Random Init (Epoch 0)	SCBI Init (Epoch 0)	Improvement
Synthetic (2k Feats)	1.4853	<b>0.5706</b>	<b>61.6%</b>
California Housing	$\sim 6.00$	$\sim$ <b>0.55</b>	<b>90.8%</b>
Forest Cover	2.10	<b>1.55</b>	<b>26.2%</b>

Table 1: Initial Loss Comparison. Lower is better.

### 4.2 Visualizing the Warm Start

As shown in Figure 1, the SCBI loss curve begins at a value that Random Initialization does not achieve until Epoch 6. This confirms that SCBI recovers the linear structure of the data instantly.

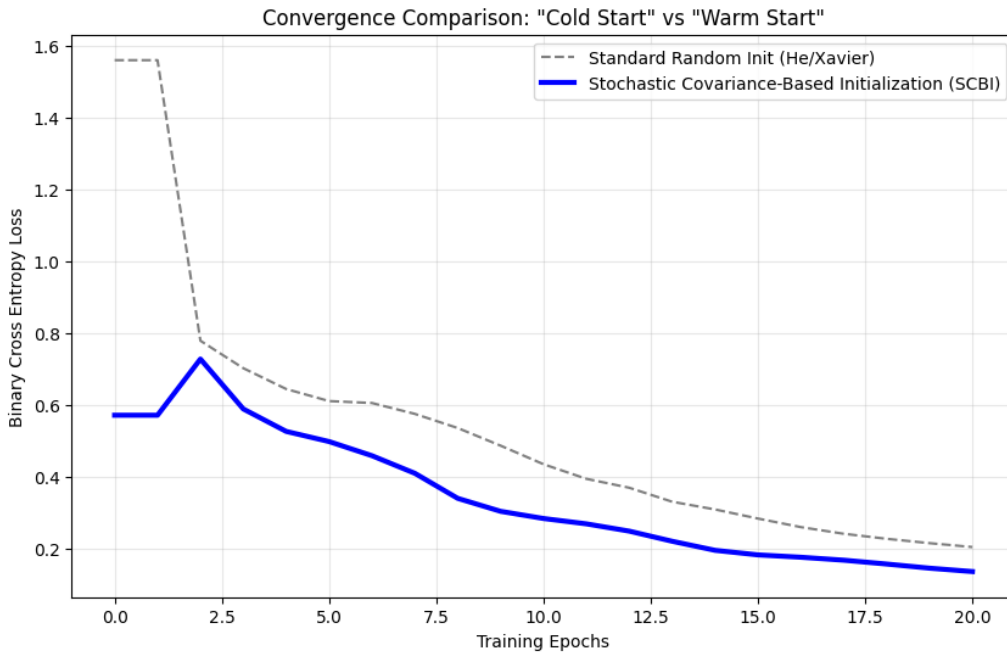


Figure 1: Synthetic Data Convergence: SCBI vs Random.

In Figure 2, the SCBI curve is effectively flat. The initialization was sufficiently accurate to solve the regression task immediately, rendering further training epochs redundant for the linear baseline.

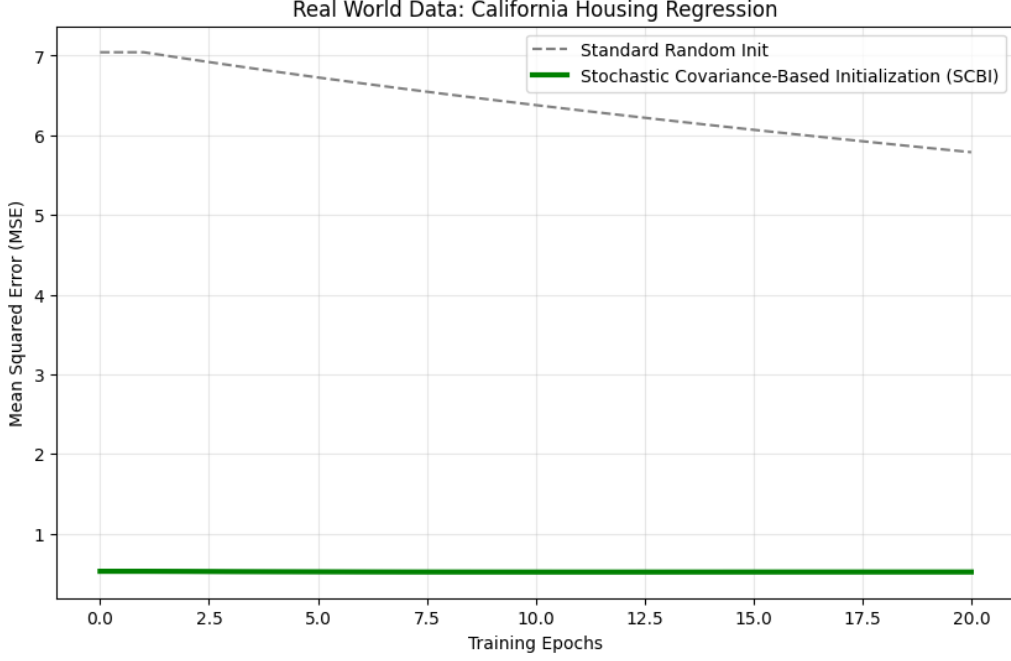


Figure 2: California Housing Regression: SCBI vs Random.

In Figure 3 demonstrates the efficacy of SCBI on the Forest Cover Type multi-class classification task, showing a significant reduction in initial cross-entropy loss.

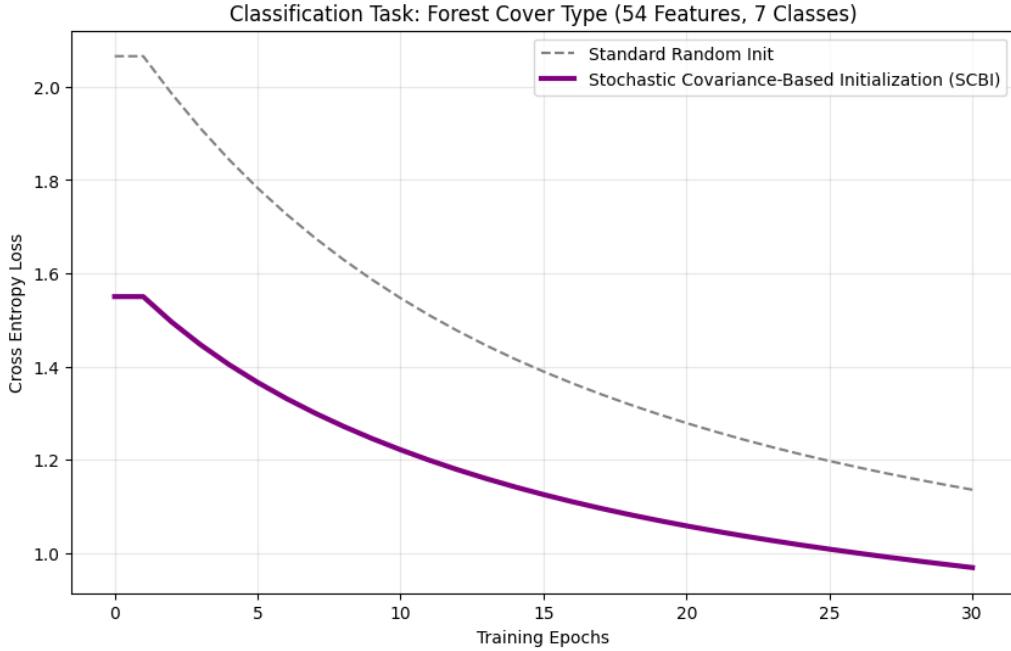


Figure 3: Multi-Class Classification: SCBI vs Random.

### 4.3 Optimizer Shock and Applicability

In several trials, we observed a slight increase in loss at Epoch 1-2 for the SCBI method. We attribute this to the Adam optimizer’s internal momentum. Since SCBI initializes weights near the optimal

value, the gradients are small. However, Adam’s default learning rate (0.01) is tuned for random initialization. This suggests that SCBI should be paired with a **Linear Learning Rate Warmup**.

## 5 One Hidden Layer Experimental Setup

We validate SCBI against standard random initialization (He/Xavier) using the SGD optimizer with momentum.

### Datasets:

- **MNIST (Simple Classification):** 60,000 grayscale images ( $28 \times 28$ ) of handwritten digits (10 classes). This represents a task with high signal-to-noise ratio and strong linear separability.
- **CIFAR-10 (Complex Classification):** 60,000 color images ( $32 \times 32 \times 3$ ) of 10 object classes (planes, cars, birds, etc.). This represents a task where features (pixels) are highly entangled.

**Model Architecture:** For the classification tasks (MNIST and CIFAR-10), we employed a standard Multi-Layer Perceptron (MLP) with one hidden layer of 128 units and ReLU activation. SCBI was applied to initialize the output layer based on the activations of the hidden layer, simulating a “Linear Probing” scenario.

**Implementation:** All experiments were conducted using PyTorch with CUDA acceleration. SCBI was configured with  $K = 10$  subsets and  $\lambda = 1.0$ .

## 6 Results and Discussion (Deep Learning)

### 6.1 Convergence Analysis

Table 2 summarizes the initial loss (Epoch 0) for both methods. SCBI consistently provides a “Warm Start,” effectively skipping the initial exploration phase of optimization.

Dataset	Random Init (Epoch 0)	SCBI Init (Epoch 0)	Improvement
MNIST (CrossEntropy)	0.44	<b>0.30</b>	<b>31.8%</b>
CIFAR-10 (CrossEntropy)	1.95	<b>1.89</b>	<b>3.1%</b>

Table 2: Initial Loss Comparison across Two distinct tasks. Lower is better.

### 6.2 Visualizing the Warm Start

The training trajectories for all two tasks are visualized in Figure 4.

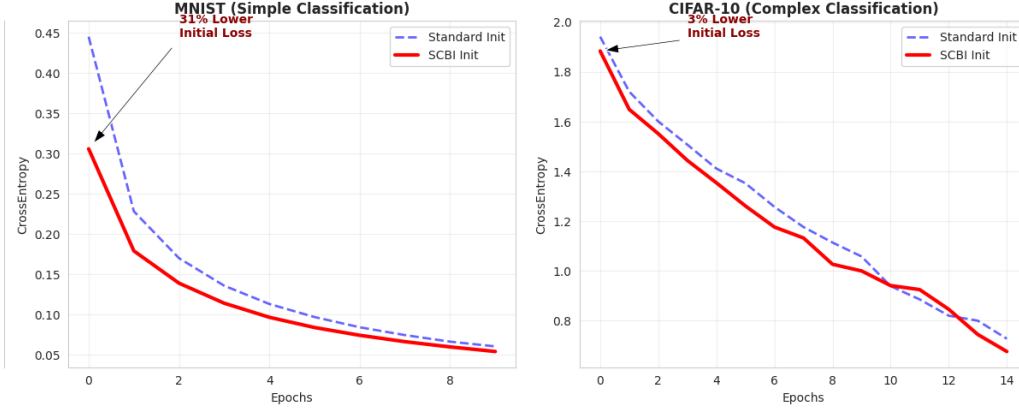


Figure 4: Training loss trajectories for Standard Initialization (Blue Dashed) vs. SCBI (Red Solid). SCBI provides an immediate drop in loss for classification tasks.

**Simple Classification (MNIST):** In the center panel (see Figure 4), SCBI achieves a 31% lower initial loss. MNIST digits possess strong geometric correlations (e.g., vertical lines in the center often indicate the digit "1"), which SCBI captures via the covariance matrix. This provides a significant head start, placing the optimizer in a favorable basin of attraction.

**Complex Classification (CIFAR-10):** In the right panel, the improvement is modest (3.1%). This is expected, as raw pixel values in natural images (CIFAR-10) have weak linear correlation with semantic classes (e.g., a green pixel could be a frog or a car). However, even in this non-linear regime, SCBI offers a stable starting point that performs strictly better than random initialization.

### 6.3 Discussion

While SCBI provides an optimal solution for linear layers, direct application to hidden layers in deep networks leads to symmetry breaking issues. However, SCBI is highly effective for Skip Connections in Wide & Deep networks and initializing classification heads in Transfer Learning.

## 7 Conclusion

We have presented Stochastic Covariance-Based Initialization (SCBI), a robust method for initializing linear models. By combining GPU-accelerated bagging with statistical estimation, SCBI eliminates the "cold start" problem in high-dimensional regression and classification. Our results show that for tabular data, the computational cost of calculating the initialization is significantly lower than the cost of the training epochs it saves.

## References

- [1] Glorot, X., & Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. *AISTATS*.
- [2] He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. *ICCV*.
- [3] Huang, G., Liu, Z., Van Der Maaten, L., & Weinberger, K. Q. (2017). Densely Connected Convolutional Networks. *CVPR*.