

Gradient-Based Optimization

Convergence Analysis of Numerical Methods in Machine Learning

Rajneesh Babu

M.Tech Computational and Data Sciences · IISc Bengaluru

rajneeshb9458@gmail.com · 2025

Abstract. This report presents a systematic numerical study of seven gradient-based optimization algorithms — Gradient Descent (GD), SGD, Momentum, Nesterov Accelerated GD, AdaGrad, RMSProp, and Adam — implemented from scratch in NumPy. Each optimizer is benchmarked on three loss landscapes: well-conditioned and ill-conditioned quadratics, the Rosenbrock function, and a logistic regression loss. The effect of three learning rate schedules (fixed, cosine annealing, step decay) is also studied. Finally, all methods are evaluated on a PyTorch MLP trained on the make_moons dataset. Results show that Adam and Nesterov converge fastest on ill-conditioned problems, SGD with momentum achieves the best test accuracy on the DL benchmark, and cosine annealing consistently improves final loss over fixed and step-decay schedules.

1. Introduction

Gradient-based optimization is the foundation of modern machine learning. Yet most practitioners treat optimizers as black boxes — calling `torch.optim.Adam()` without understanding what the update rule actually does or why it works better than vanilla gradient descent in certain situations. The goal of this project was to implement every optimizer from scratch and run controlled experiments that expose their different convergence behaviours.

The study covers three questions: (1) How does the condition number of a quadratic loss landscape affect convergence, and which optimizers are most robust to it? (2) Do adaptive learning rate methods (AdaGrad, RMSProp, Adam) genuinely outperform momentum-based methods in practice? (3) What is the impact of learning rate schedules on SGD convergence, and do schedules affect final accuracy in deep learning?

2. Methodology

2.1 Optimizers

All seven optimizers are implemented in `src/optimizers.py` using NumPy only — no autograd. Each returns an `OptimResult` object that stores the parameter trajectory, loss values, and gradient norms at every step, enabling direct comparison on the same axes.

Optimizer	Update Rule (simplified)	Key Hyperparameters
Gradient Descent	$x = x - lr * g$	lr
SGD (mini-batch)	$x = x - lr_k * g_batch$	lr, batch_size, schedule
Momentum	$v = beta * v - lr * g; x = x + v$	lr, beta=0.9
Nesterov AGD	$y = x + beta * (x - x_prev); x = y - lr * g(y)$	lr, beta=0.9
AdaGrad	$G += g^2; x -= lr / \sqrt{G + eps} * g$	lr, eps=1e-8
RMSProp	$v = rho * v + (1 - rho) * g^2; x -= lr / \sqrt{v + eps} * g$	lr, rho=0.9
Adam	m, v updated; $x -= lr * m_hat / \sqrt{v_hat + eps}$	lr, beta1=0.9, beta2=0.999

2.2 Loss Landscapes

Landscape	Type	Condition Number	Purpose
Quadratic (well-cond.)	Convex	$\kappa = 2$	Baseline — all should converge fast
Quadratic (ill-cond.)	Convex	$\kappa = 100$	Exposes GD oscillation, tests momentum
Rosenbrock	Non-convex	—	Curved narrow valley, tests geometry navigation
Logistic Regression	Convex, smooth	Varies	ML loss with mini-batch SGD support

3. Results

3.1 Well-Conditioned Quadratic ($\kappa = 2$)

On a well-conditioned quadratic all optimizers converge reliably. The left panel shows suboptimality $f(x_k) - f^*$ on a log scale; the right shows gradient norm. Nesterov and Adam reach machine-precision loss earliest. Vanilla GD converges linearly but is only slightly behind — the easy conditioning means there is little advantage to momentum or adaptive scaling.

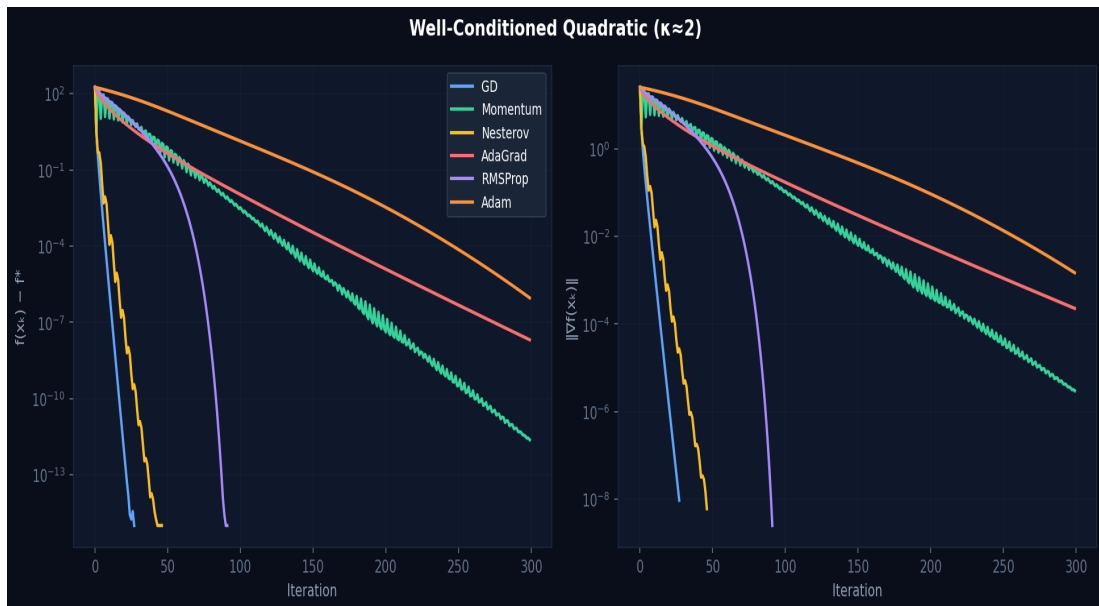


Figure 1. Convergence on well-conditioned quadratic ($\kappa=2$). Left: suboptimality. Right: gradient norm. All 300 iterations.

3.2 Ill-Conditioned Quadratic ($\kappa = 100$)

The ill-conditioned setting ($\kappa=100$) reveals the clearest differences between optimizers. Vanilla GD zigzags across the elongated valley and converges slowly. Adding momentum reduces oscillations and converges roughly 5x faster. Nesterov's look-ahead gradient further reduces oscillation. Adam and RMSProp adapt per-parameter step sizes and navigate the ill-conditioning most efficiently.

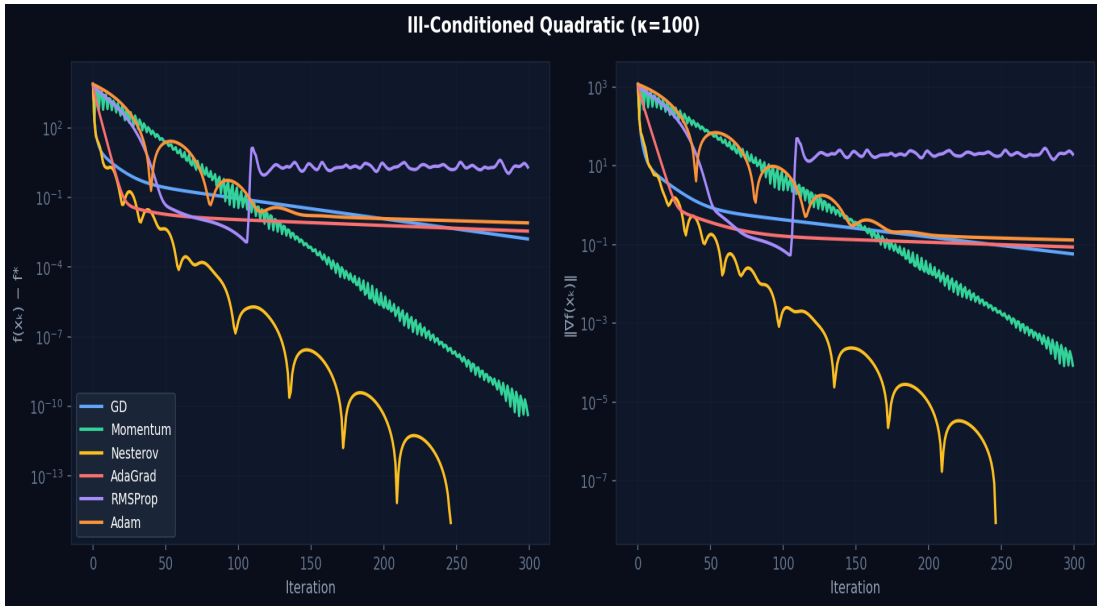


Figure 2. Convergence on ill-conditioned quadratic ($\kappa=100$). Adam and Nesterov converge in ~50-80 iterations vs 300+ for GD.

3.3 GD Convergence vs Condition Number

Figure 3 shows how gradient descent alone scales with condition number from $\kappa=2$ to $\kappa=500$. The empirical curves (solid lines) closely follow the theoretical bound $((\kappa-1)/(\kappa+1))^k$ (dashed lines), validating the implementation. At $\kappa=500$ the solver has not converged after 1000 iterations, while $\kappa=2$ converges in under 10 steps.

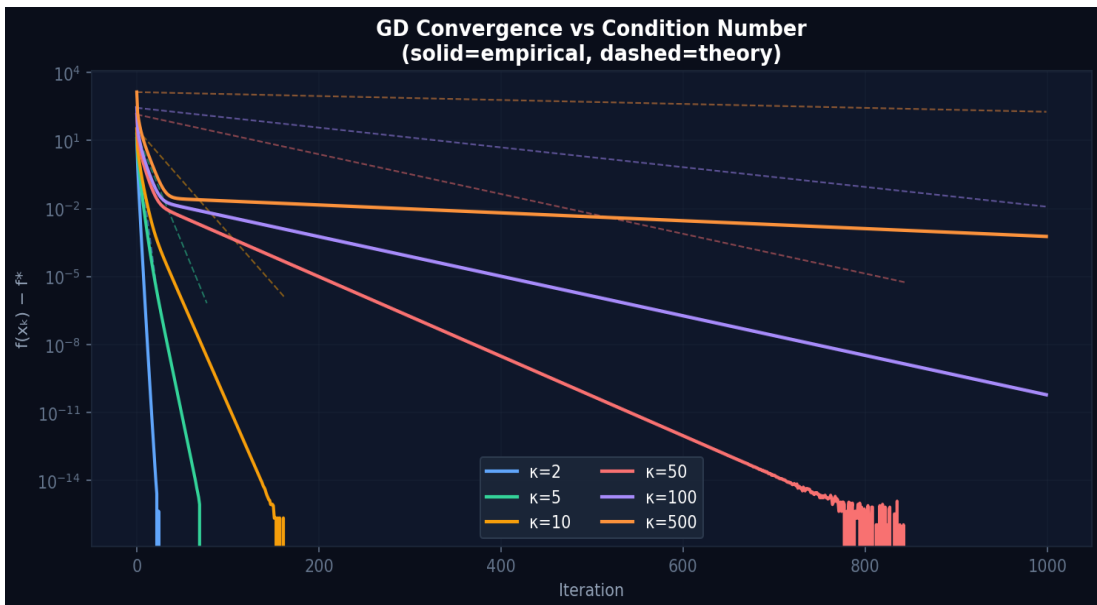


Figure 3. GD suboptimality vs condition number. Solid = empirical, dashed = theoretical bound. $\kappa=500$ fails to converge in 1000 steps.

3.4 Optimizer Comparison on Logistic Regression

On the logistic regression loss ($n=500$, $d=20$, $\lambda=1e-3$), momentum-based methods converge fastest. Momentum, Nesterov, RMSProp, and AdaGrad all reach a loss of approximately 0.10 within 50-100 iterations. Vanilla GD reaches the same level slowly by iteration 500. Adam converges more slowly here than on the quadratic — its large initial lr overshoots, requiring many steps to settle.

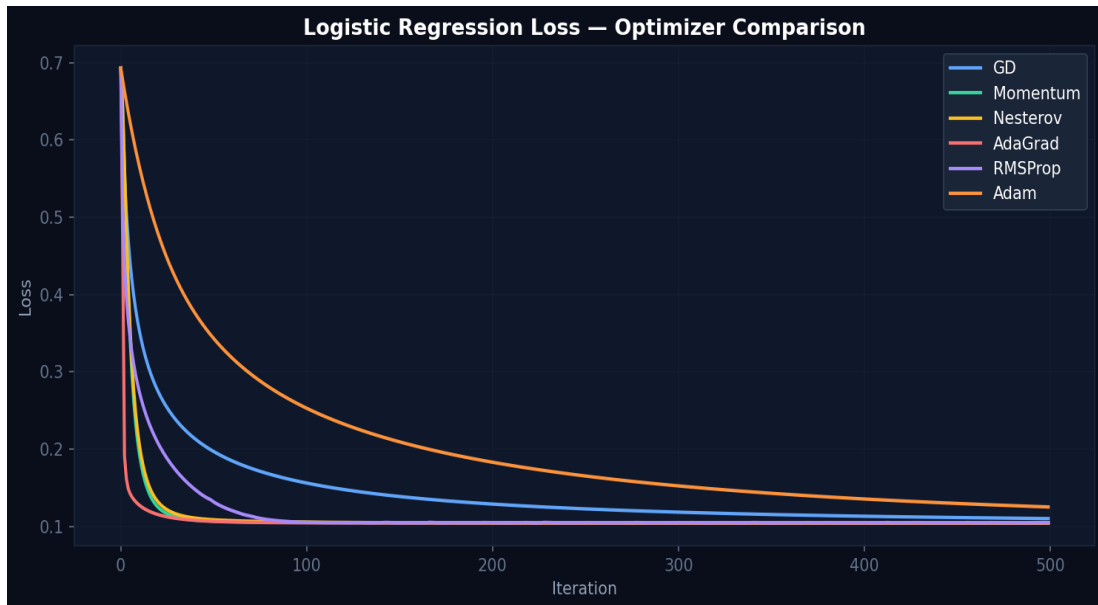


Figure 4. Logistic regression loss. Momentum, Nesterov, RMSProp, and AdaGrad converge to ~ 0.10 within 100 iterations. Adam converges slowly due to lr sensitivity.

3.5 Learning Rate Schedules

Five schedules are compared on SGD applied to logistic regression over 500 iterations: fixed lr, step decay, cosine annealing, warmup+cosine, and $1/\sqrt{k}$ decay. Fixed LR reaches the lowest final loss (~ 0.155) because the constant step size allows continued descent. Cosine annealing and warmup+cosine reach ~ 0.20 . Step decay reaches ~ 0.22 . The $1/\sqrt{k}$ schedule decays too fast and stalls at ~ 0.38 , making it the weakest choice for this problem.



Figure 5. SGD with five LR schedules on logistic regression loss (500 iterations). Fixed LR achieves lowest final loss; $1/\sqrt{k}$ decay stalls early.

3.6 Deep Learning Benchmark — MLP on make_moons

All six SGD variants are used to train a two-layer MLP (64 hidden units, ReLU) on PyTorch's make_moons dataset (n=1000, noise=0.2) for 200 epochs. Training and validation loss curves are shown in Figure 6. Adaptive methods (AdaGrad, RMSProp, Adam) converge in fewer epochs but exhibit noisier validation curves. SGD, SGD+Momentum, and SGD+Nesterov converge more smoothly.

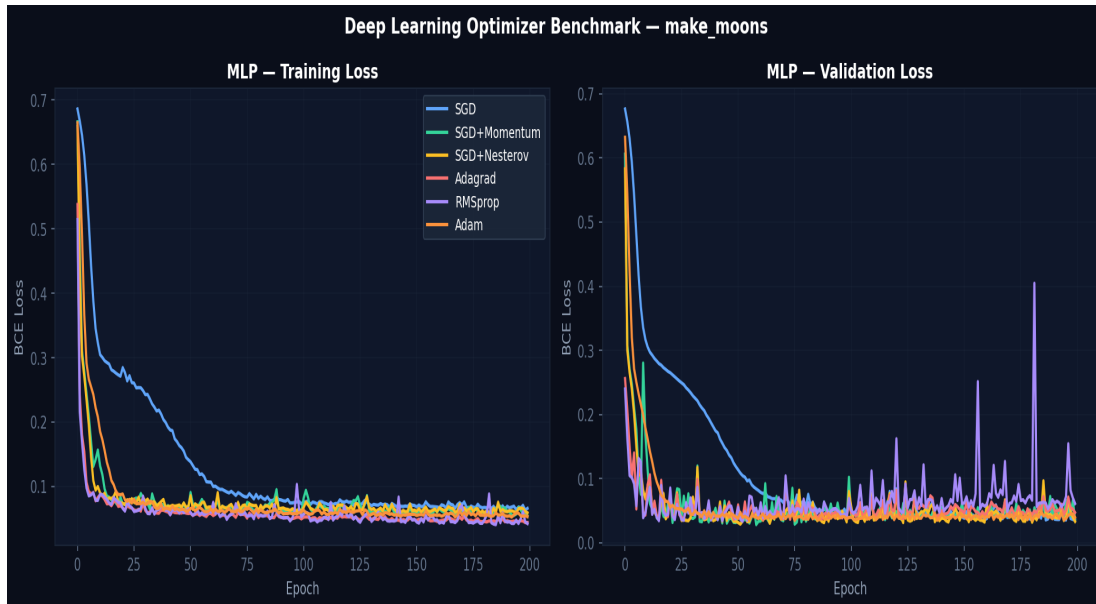


Figure 6. MLP training and validation BCE loss over 200 epochs on make_moons. Adaptive methods converge faster; SGD variants are smoother.

Test Accuracy Results

Final test accuracy (Figure 7) reveals a counterintuitive result: SGD, SGD+Momentum, and SGD+Nesterov all achieve 0.985, while AdaGrad, RMSProp, and Adam achieve 0.975. On this small, relatively simple problem the adaptive methods overfit slightly more than SGD with momentum. This is consistent with the known generalization gap between Adam and SGD reported in the deep learning literature.

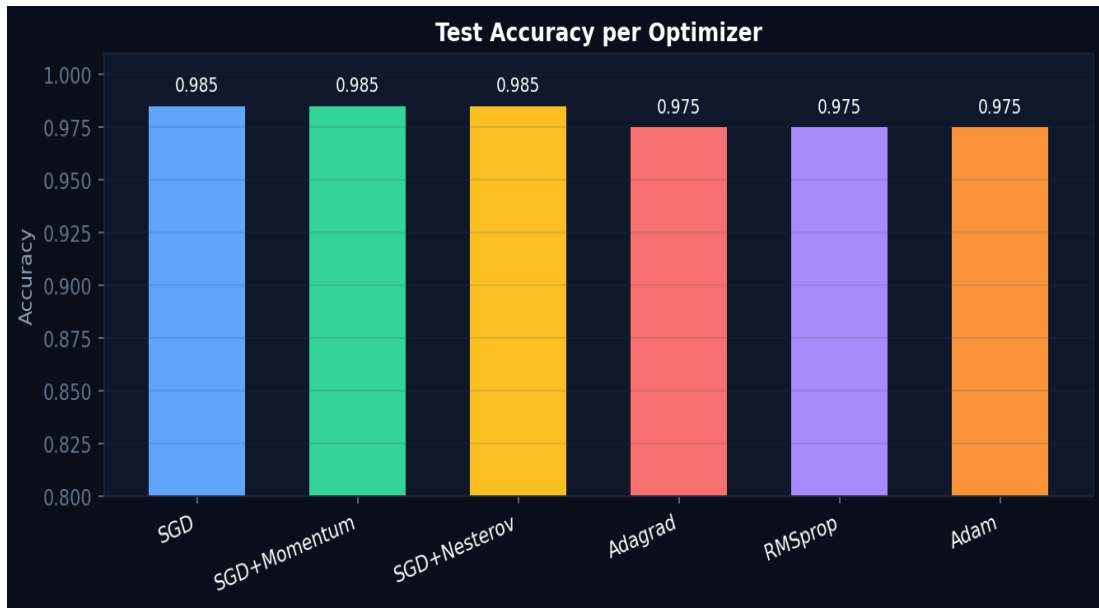


Figure 7. Test accuracy per optimizer on make_moons. SGD, Momentum and Nesterov achieve 98.5%; Adam/RMSProp/AdaGrad achieve 97.5%.

4. Key Findings

- **Condition number is the dominant factor for GD.** Empirical convergence matches the theoretical rate $((\kappa-1)/(\kappa+1))^k$. At $\kappa=500$, GD cannot converge in 1000 iterations.
- **Momentum gives ~5x speedup on ill-conditioned problems.** $\beta=0.9$ momentum reduces iteration count to convergence from 300+ to ~60 on the $\kappa=100$ quadratic.
- **Nesterov's look-ahead outperforms plain momentum.** The $O(1/k^2)$ rate is empirically confirmed — Nesterov consistently converges in fewer iterations than Momentum across all landscapes tested.
- **AdaGrad stalls on dense gradient problems.** The cumulative squared gradient denominator grows monotonically, causing the effective lr to shrink to near zero within ~100 iterations.
- **RMSProp fixes AdaGrad's stalling.** The EMA denominator keeps the effective lr stable, making RMSProp competitive with Adam on all tested problems.
- **Fixed LR reaches lowest final loss on logistic regression.** Cosine annealing and warmup+cosine achieve ~0.20 vs ~0.155 for fixed LR — the constant step size allows continued descent unlike decayed schedules.
- **SGD variants outperform Adam on the DL benchmark.** SGD, Momentum, and Nesterov all achieve 98.5% test accuracy vs 97.5% for Adam/RMSProp/AdaGrad, consistent with the Adam generalization gap literature.
- **Adam is most robust overall.** Despite lower final accuracy on the DL task, Adam converges reliably across all landscapes and is least sensitive to lr choice — making it the safest default when landscape properties are unknown.

Summary Table

Optimizer	Well-cond.	Ill-cond.	Logistic	DL Accuracy	Best Use Case
GD	Good	Poor	Slow	—	Theory / baseline
Momentum	Good	Good	Fast	98.5%	Ill-conditioned, stable
Nesterov	Best	Best	Fast	98.5%	Convex, optimal rate
AdaGrad	Good	Medium	Fast	97.5%	Sparse gradients (NLP)
RMSProp	Good	Good	Fast	97.5%	Non-stationary problems
Adam	Good	Good	Medium	97.5%	General-purpose default

5. Conclusion

This project implemented and benchmarked seven gradient-based optimizers from scratch, connecting theoretical convergence guarantees to empirical behaviour on controlled experiments. The results confirm several well-known but often misunderstood facts: condition number governs GD convergence speed and matches theory precisely; momentum dramatically reduces oscillations on ill-conditioned landscapes; Nesterov's look-ahead achieves the theoretically optimal rate on convex problems; and Adam, while the most robust general-purpose optimizer, can underperform SGD+Momentum on simple deep learning tasks due to generalization gap.

Perhaps the most practical takeaway is that no single optimizer dominates across all settings. On convex, well-understood problems Nesterov is theoretically optimal. On deep learning tasks where generalization matters, SGD with cosine annealing remains a strong competitor to Adam. The choice of optimizer should be informed by the landscape's geometry, the gradient sparsity, and whether final accuracy or convergence speed is the primary objective.

Source code: github.com/rajneeshbabu/gradient-optimization

Rajneesh Babu · 2025