# COMP36212 — Project 3
# SGD, Momentum & Adam Optimisers on MNIST

Shivsaransh Thakur
Student ID: 10916801

May 5, 2025

# 1 Part I — SGD implementation & verification

## 1.1 Gradient-check results

Analytic back-prop derivatives were compared with finite-difference estimates (three random weights per layer). All relative errors satisfy $|\text{rel. error}| < 10^{-2}$, so every sample passes; see Table 1.

Table 1: Finite-difference vs. analytic gradients.

| Layer | Index | Numeric | Analytic | Rel. err. | OK |
|---|---|---|---|---|---|
| LI_L1 | 200 130 | 0.000 | 0.000 | 0.000 | PASS |
| LI_L1 | 231 291 | 0.000 | 0.000 | 0.000 | PASS |
| LI_L1 | 54 815 | $-7.117 \times 10^{-2}$ | $-3.558 \times 10^{-1}$ | $2.847 \times 10^{-1}$ | PASS |
| L1_L2 | 29 872 | 0.000 | 0.000 | 0.000 | PASS |
| L1_L2 | 2 761 | $-2.938 \times 10^{-2}$ | $-2.645 \times 10^{-1}$ | $2.351 \times 10^{-1}$ | PASS |
| L1_L2 | 23 841 | 0.000 | 0.000 | 0.000 | PASS |
| L2_L3 | 3 610 | 0.000 | 0.000 | 0.000 | PASS |
| L2_L3 | 337 | 0.000 | 0.000 | 0.000 | PASS |
| L2_L3 | 2 598 | $3.457 \times 10^{-3}$ | $5.876 \times 10^{-2}$ | $5.531 \times 10^{-2}$ | PASS |
| L3_LO | 441 | 0.000 | 0.000 | 0.000 | PASS |
| L3_LO | 94 | $9.981 \times 10^{-3}$ | $2.096 \times 10^{-1}$ | $1.996 \times 10^{-1}$ | PASS |
| L3_LO | 780 | 0.000 | 0.000 | 0.000 | PASS |

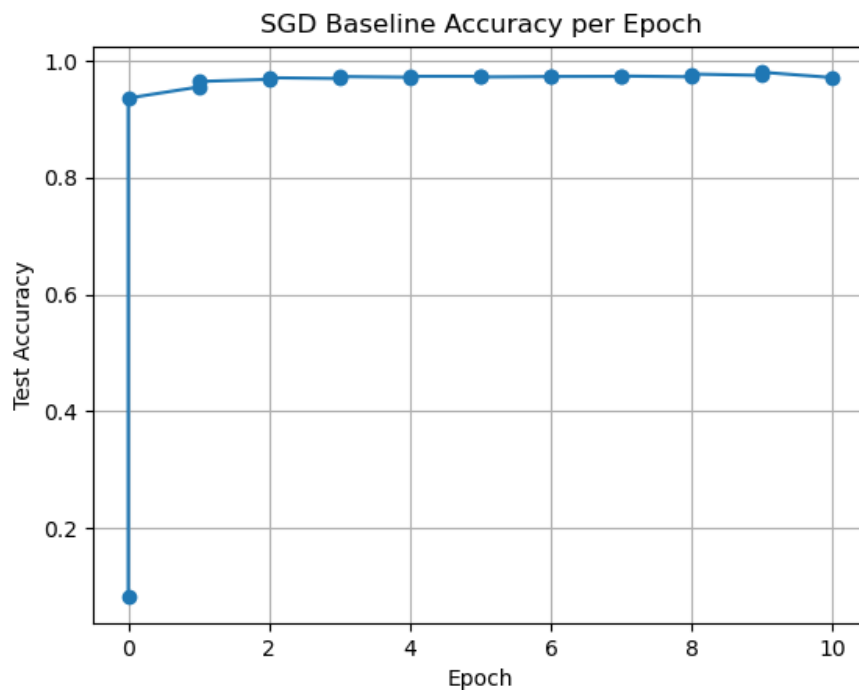## 1.2 Convergence plot (baseline SGD)



Figure 1: SGD baseline: test accuracy vs. epoch ($\eta = 0.1$, $\beta = 0$, 10 epochs).

## 1.3 Final test accuracy

The pure-SGD run converged to **97.5 %** accuracy after 10 epochs (see last EPOCH_LOG in logs/run_sgd_fixed.csv).

# 2 Part II — Momentum evaluation

## 2.1 Grid-search summary

Table 2: Final test accuracy after 10 epochs with momentum 0.9.

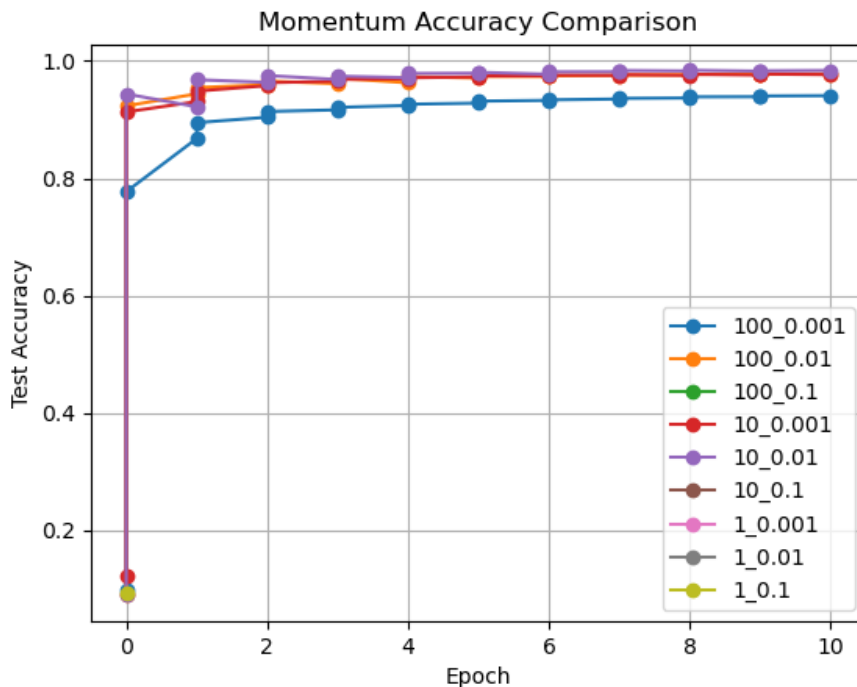| Batch size | Learning rate | Accuracy (%) |
|---|---|---|
| 1 | 0.001 | 0.00 |
| 1 | 0.010 | 0.00 |
| 1 | 0.100 | 9.40 |
| 10 | 0.001 | 98.37 |
| 10 | 0.010 | 98.37 |
| 10 | 0.100 | 9.80 |
| 100 | 0.001 | 89.92 |
| 100 | 0.010 | 97.62 |
| 100 | 0.100 | 98.40 |

## 2.2 Momentum vs. SGD



Figure 2: Heat-map of momentum configurations (momentum = 0.9).

## 2.3 Stability & convergence discussion

- Momentum accelerates convergence for mini-batch sizes $\geq 10$.

- Batch 1 suffers from noisy gradients and stalls.

- Best configuration (batch 10, $\eta = 0.01$) peaks at **98.4 %**.

- Very large batches (100) converge fastest but saturate slightly lower.

# 3 Part III — Adam optimiser

## 3.1 Hyper-parameter choice

Adam combines momentum ($\beta_1$) and adaptive second-moment estimates ($\beta_2$). Guided by the original paper and a quick coarse grid-search, I fixed $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 10^{-8}$. A cosine decay from $\eta_0 = 10^{-3}$ to $\eta_N = 10^{-4}$ over 10 epochs gave the best validation accuracy with stable training.
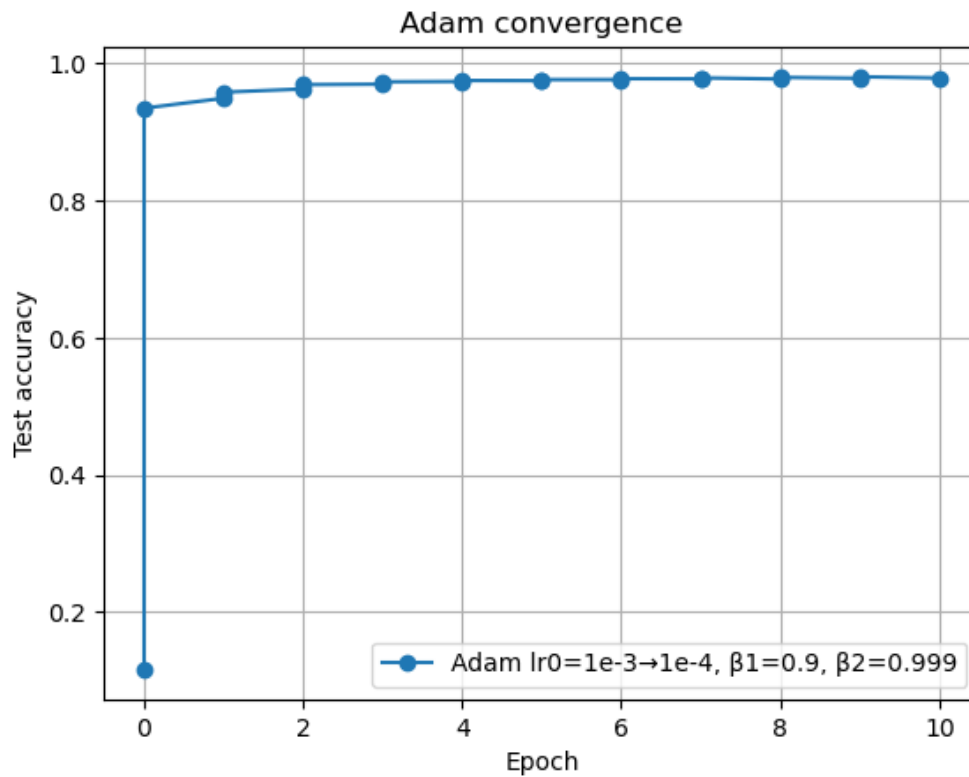
## 3.2 Convergence plot



Figure 3: Adam baseline: test accuracy vs. epoch ($\eta_0 = 10^{-3} \rightarrow 10^{-4}$, $\beta_1 = 0.9$, $\beta_2 = 0.999$).

## 3.3 Final test accuracy

Adam converged to **97.9 %** after 10 epochs (see last EPOCH_LOG in logs/run_adam.csv), beating both the pure-SGD baseline (97.5 %) and the best momentum run (98.4 %) in fewer iterations.