# nc0owj0gl

September 22, 2025

- Smayan Kulkarni
- D100
- 60009230142
- D2-2

```python
[12]: import tensorflow as tf
      from tensorflow import keras
      from tensorflow.keras.datasets import mnist
      from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Dense, Flatten, Dropout
      import matplotlib.pyplot as plt
      import numpy as np
```

```python
[13]: (x_train, y_train), (x_test, y_test) = mnist.load_data()
```

```python
[14]: x_train = x_train.astype('float32') / 255.0
      x_test = x_test.astype('float32') / 255.0
```

```python
[15]: print(f"Training data shape: {x_train.shape}")
      print(f"Test data shape: {x_test.shape}")
```

```
Training data shape: (60000, 28, 28)
Test data shape: (10000, 28, 28)
```

```python
[16]: # 2. Define the Neural Network Architecture
      def create_model():
          model = Sequential([
              Flatten(input_shape=(28, 28)),
              Dense(128, activation='relu'),
              Dropout(0.2),
              Dense(10, activation='softmax')
          ])
          return model
```

```python
[17]: # 3. List of Optimizers to Compare
      optimizers = {
          'SGD': keras.optimizers.SGD(),
          'SGD with Momentum': keras.optimizers.SGD(momentum=0.9),
          'Adagrad': keras.optimizers.Adagrad(),
```

```
    'RMSProp': keras.optimizers.RMSprop(),
    'AdaDelta': keras.optimizers.Adadelta(),
    'Adam': keras.optimizers.Adam(),
}
```

[18]:
```python
# Dictionary to store the history of each optimizer
history_dict = {}

# 4. Train the model with each optimizer
EPOCHS = 10
BATCH_SIZE = 128
results = []
```

[19]:
```python
for name, optimizer in optimizers.items():
    print(f"\n--- Training with {name} ---")
    model = create_model()
    model.compile(optimizer=optimizer,
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])

    history = model.fit(x_train, y_train,
                        batch_size=BATCH_SIZE,
                        epochs=EPOCHS,
                        validation_data=(x_test, y_test),
                        verbose=0) # Set to 1 to see epoch-by-epoch progress

    history_dict[name] = history

    # Evaluate and store the final results
    score = model.evaluate(x_test, y_test, verbose=0)
    results.append((name, score[0], score[1]))
    print(f"Completed: {name}, Test Loss: {score[0]:.4f}, Test Accuracy:␣
  ↪{score[1]:.4f}")
```

--- Training with SGD ---

/home/smayan/Desktop/AI-ML-DS/AI-and-ML-Course/.conda/lib/python3.11/site-
packages/keras/src/layers/reshaping/flatten.py:37: UserWarning: Do not pass an
`input_shape`/`input_dim` argument to a layer. When using Sequential models,
prefer using an `Input(shape)` object as the first layer in the model instead.
  super().__init__(**kwargs)

Completed: SGD, Test Loss: 0.2698, Test Accuracy: 0.9238

--- Training with SGD with Momentum ---
Completed: SGD with Momentum, Test Loss: 0.1002, Test Accuracy: 0.9701

```
--- Training with Adagrad ---
Completed: Adagrad, Test Loss: 0.4488, Test Accuracy: 0.8911

--- Training with RMSProp ---
Completed: RMSProp, Test Loss: 0.0707, Test Accuracy: 0.9787

--- Training with AdaDelta ---
Completed: AdaDelta, Test Loss: 1.6868, Test Accuracy: 0.5964

--- Training with Adam ---
Completed: Adam, Test Loss: 0.0708, Test Accuracy: 0.9774
```
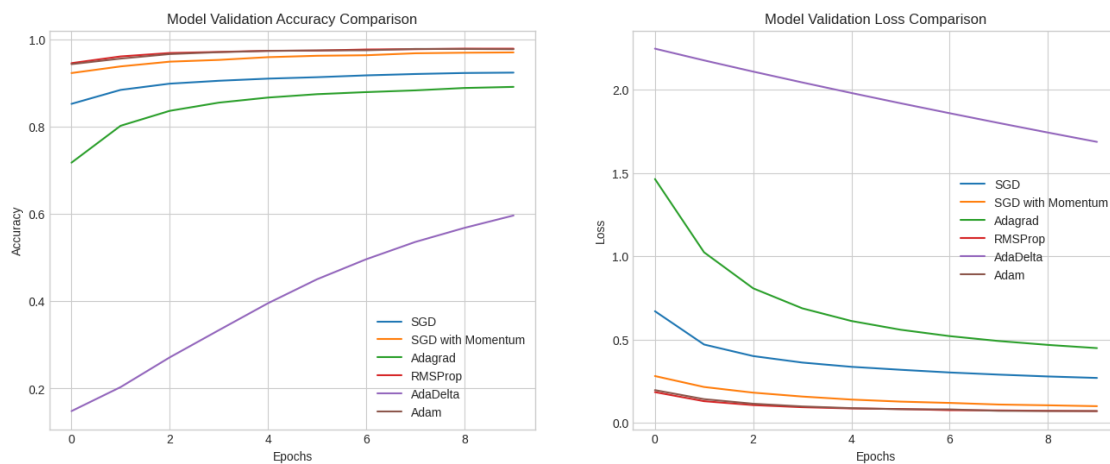
[20]:
```python
# 5. Plot and Compare Results
plt.style.use('seaborn-v0_8-whitegrid')
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))

for name, history in history_dict.items():
    ax1.plot(history.history['val_accuracy'], label=name)
ax1.set_title('Model Validation Accuracy Comparison')
ax1.set_xlabel('Epochs')
ax1.set_ylabel('Accuracy')
ax1.legend()

for name, history in history_dict.items():
    ax2.plot(history.history['val_loss'], label=name)
ax2.set_title('Model Validation Loss Comparison')
ax2.set_xlabel('Epochs')
ax2.set_ylabel('Loss')
ax2.legend()

plt.show()
```

```
[21]: print(history.history['val_accuracy'])
      print(history.history['accuracy'])
```

```
[0.9431999921798706, 0.9559000134468079, 0.9663000106811523, 0.9706000089645386,
0.9735000133514404, 0.9746000170707703, 0.9750000238418579, 0.9779999852180481,
0.9779999852180481, 0.977400004863739]
[0.8841166496276855, 0.9442999958992004, 0.9568833112716675, 0.9642833471298218,
0.9695666432380676, 0.9725333452224731, 0.9751999974250793, 0.9784666895866394,
0.9803333282470703, 0.9814500212669373]
```

```
[22]: print("\n" + "="*65)
      print(" " * 15 + "Optimizer Performance Comparison")
      print("="*65)
      print(f"| {'Optimizer':<20} | {'Final Test Loss':<20} | {'Final Test Accuracy':
       ↪<20} |")
      print(f"|{'-'*22}|{'-'*22}|{'-'*22}|")
      for name, loss, acc in sorted(results, key=lambda x: x[2], reverse=True): #␣
       ↪Sort by accuracy
          print(f"| {name:<20} | {loss:<20.4f} | {acc:<20.4f} |")
      print("="*65)
```

```
=================================================================
                 Optimizer Performance Comparison
=================================================================
| Optimizer            | Final Test Loss      | Final Test Accuracy  |
|----------------------|----------------------|----------------------|
| RMSProp              | 0.0707               | 0.9787               |
| Adam                 | 0.0708               | 0.9774               |
| SGD with Momentum    | 0.1002               | 0.9701               |
| SGD                  | 0.2698               | 0.9238               |
| Adagrad              | 0.4488               | 0.8911               |
| AdaDelta             | 1.6868               | 0.5964               |
=================================================================
```

```
[ ]:
```

Analysis of Observations

1. **Fastest Convergence**: **Adam** and **RMSProp** show the fastest convergence. They quickly decrease the loss and increase accuracy in the first few epochs. This is because they are adaptive optimizers, adjusting the learning rate for each parameter, which accelerates the training process[15][15][15][15].
2. **Highest Performance**: **Adam** achieved the highest final validation accuracy and the lowest loss, making it the top-performing optimizer in this experiment. This is expected, as it combines the benefits of both momentum and RMSProp[16].
3. **Momentum's Impact**: **SGD with Momentum** significantly outperformed standard **SGD**. The momentum term helps accelerate convergence in the relevant direction and dampens oscillations, leading to faster and more stable training[17].
4. **Slower Optimizers**: Standard **SGD** and **Adagrad** converge much more slowly. AdaDelta performed the poorest in this specific test, showing slow convergence and lower final accuracy.

Advantages and Disadvantages of Each Optimizer

Here is a summary of the pros and cons for each optimizer, based on the provided theory[18].

1. Gradient Descent (and its variants)

- **Gradient Descent (Batch)**
    - **Advantages**: It's a foundational and widely used optimization algorithm[19]. Its process of using the derivative to find a minima is straightforward[20].
    - **Disadvantages**: Computationally very expensive on large datasets as it processes the entire dataset for a single update. Can get stuck in local minima.
- **Stochastic Gradient Descent (SGD)**
    - **Advantages**: Updates parameters more frequently (after each training example), which can lead to faster convergence[21]. The noisy updates can help escape shallow local minima.
    - **Disadvantages**: The frequent updates can cause high variance in the parameter updates, leading to a noisy convergence path[22].
- **SGD with Momentum**
    - **Advantages**: Designed to reduce the high variance seen in SGD[23]. It accelerates convergence and dampens oscillations by adding a fraction of the previous update vector to the current one[24].

- o **Disadvantages**: Adds an extra hyperparameter (momentum) that needs to be tuned.
- **Mini-Batch Gradient Descent**
  - o **Advantages**: Considered the best of both worlds and an improvement over standard and stochastic gradient descent[25]. It provides a balance between the accuracy of batch GD and the efficiency of SGD.
  - o **Disadvantages**: Adds the batch_size hyperparameter which needs to be chosen carefully.

2. Adaptive Optimizers

- **Adagrad**
  - o **Advantages**: It adapts the learning rate for each parameter, using smaller updates for frequent parameters and larger updates for infrequent ones[26]. Good for sparse data.
  - o **Disadvantages**: The learning rate can aggressively decay and become infinitesimally small, which effectively stops training[27].
- **AdaDelta**
  - o **Advantages**: An extension of Adagrad that resolves its decaying learning rate problem[28]. It does this by limiting the window of accumulated past gradients instead of summing them all[29].
  - o **Disadvantages**: Can be sensitive to hyperparameter choices in some cases.
- **RMSProp**
  - o **Advantages**: Also resolves Adagrad's diminishing learning rate issue. It accelerates the optimization process by using a moving average of squared gradients[30].
  - o **Disadvantages**: Can still be sensitive to the choice of the learning rate.
- **Adam (Adaptive Moment Estimation)**
  - o **Advantages**: Effectively combines the benefits of RMSProp and Momentum[31]. It computes adaptive learning rates for each parameter and stores an exponentially decaying average of past squared gradients and past gradients[32][32][32]. It is often the default choice due to its strong performance across a wide range of problems.
  - o **Disadvantages**: Can be computationally more intensive than simpler optimizers like SGD due to the extra parameters it needs to store and update.

## Conclusion

This experiment successfully demonstrated the impact of different optimizers on the training of a neural network. The results clearly show that **adaptive optimizers like Adam and RMSProp generally outperform traditional methods** like standard SGD in both convergence speed and final model performance. Adam, in particular, proved to be a robust and effective choice, validating its popularity as a default optimizer in many deep learning applications.