# Evaluating Optimization Methods for Predicting Football Outcomes with Feedforward Neural Networks

Julius Graf
julius.graf@berkeley.edu

Louis Sallé-Tourne
louis_salle-tourne@berkeley.edu

December 15, 2024

**Abstract**

This academic project, inspired by the "Challenge Data - Football: Can you guess the winner? by QRT" competition, aims to explore the behavior of neural network learning processes, with a particular focus on the optimization algorithms used during training. The project is part of the UC Berkeley Fall 2024 graduate course INDENG 240: Optimization Analytics, taught by Prof. Phillip Kerger. By analyzing a dataset containing various information about football teams and their match outcomes, we will train neural networks to predict the winners of football matches. This will allow us to apply and evaluate the gradient descent techniques learned in class.

## Contents

## 1 Introduction

We aim to contribute to the 2024 "Challenge Data - Football: Can you guess the winner? by QRT"Football: Can you guess the winner? by QRT. Given that this project is part of the INDENG 240 course, our primary focus will be on the optimization aspects: defining the optimization problem, modeling through neural networks, and selecting appropriate optimization algorithms.

The competition organizers provide the training dataset for this challenge directly. Two tables give the input data:

- `train_home_team_statistics_df`: This table contains all information about the home team, with `game_id` as the key. It includes a variety of statistical data such as shots taken, passes made, and other performance metrics,

- `train_away_team_statistics_df`: This table has the same structure as the home team table but contains data for the away team, also keyed by `game_id`.

These tables can be joined using the `game_id` to create a comprehensive dataset for each match. The entire dataset comprises 12,203 rows, with 142 columns for each team (home and away). Given that we split the training data into training, testing, and validation data with an 80% ratio, the training set denoted $\mathcal{D}_m$ is of size $m = 7873$. The columns are predominantly numerical, except for the first two: `LEAGUE` and `TEAM_NAME`. Thus, the input size is $p = 280$ (dropping the league and the team name, one has 140 regressors for each of both teams). Examples of numerical columns include:

- `TEAM_SHOTS_TOTAL_season_sum`: Total shots taken by the team throughout the whole season,

- `TEAM_SHOTS_INSIDEBOX_season_sum`: Total shots taken inside the box by the team throughout the whole season,

- `TEAM_SHOTS_OFF_TARGET_season_sum`: Total off-target shots taken by the team throughout the whole season,

- `TEAM_SHOTS_ON_TARGET_season_sum`: Total on-target shots taken by the team throughout the whole season,

- `TEAM_SHOTS_OUTSIDEBOX_season_sum`: Total shots taken outside the box by the team throughout the whole season,

- `TEAM_PASSES_season_sum`: Total number of passes made by the team throughout the whole season,

- `TEAM_SUCCESSFUL_PASSES_season_sum`: Total number of successful passes made by the team throughout the whole season,

The label dataset, `Y_train`, contains the same number of rows (12,203) and includes the columns `HOME_WINS`, `DRAW`, and `AWAY_WINS`, with a value of 1 indicating the match result. Therefore, the whole dataset size is 12,203 Our task is a multiclass classification problem (thus, the number of classes is $c = 3$), where we aim to predict the outcome of each football match.

To achieve this, we will utilize the `PyTorch` library to build and train our neural network. This project will allow us to apply and evaluate the gradient descent techniques and other optimization algorithms learned in class. By focusing on the optimization aspects, we hope to gain insights into how different optimization strategies affect the performance of neural networks in predicting football match outcomes.

# 2 Mathematical formalization

## 2.1 Optimization problem

We consider a neural network $f_\theta \colon x \in \mathbb{R}^p \mapsto W_2 \sigma(W_1 x + b_1) + b_2 \in \mathbb{R}^c$, where the learnable parameters are $\theta = (W_1, W_2, b_1, b_2)$. These parameters belong to the space

$$\Theta = \mathcal{M}_{h,p}(\mathbb{R}) \times \mathcal{M}_{c,h}(\mathbb{R}) \times \mathbb{R}^h \times \mathbb{R}^c.$$

Here, $W_1 \in \mathcal{M}_{h,p}(\mathbb{R})$ and $W_2 \in \mathcal{M}_{c,h}(\mathbb{R})$ are weight matrices, while $b_1 \in \mathbb{R}^h$ and $b_2 \in \mathbb{R}^c$ are bias vectors. The activation function is set to $\sigma = \mathrm{ReLU}$. For practical purposes, the parameter space $\Theta$ can be assimilated into $\mathbb{R}^d$ with

$$d = h(1 + p + c) + c.$$

This transformation is achieved by flattening the matrices $W_1$ and $W_2$ using a column-wise vectorization function $\phi_{p,q} \colon \mathcal{M}_{p,q}(\mathbb{R}) \to \mathbb{R}^{pq}$, defined as

$$\phi_{p,q}(A) = \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_q \end{bmatrix},$$

where $A_i \in \mathbb{R}^p$ is the $i$-th column of $A \in \mathcal{M}_{p,q}(\mathbb{R})$. Thus, $W_1 \equiv \phi_{h,p}(W_1)$ and $W_2 \equiv \phi_{c,h}(W_2)$. We adopt the cross-entropy loss function, defined as follows:

$$\forall \mathcal{D} \subseteq \mathbb{R}^p \times \mathbb{R}^c, |\mathcal{D}| < \infty, \quad \forall \theta \in \mathbb{R}^d, \quad l(\theta, \mathcal{D}) = -\frac{1}{m} \sum_{(x,y) \in \mathcal{D}} \sum_{j=1}^{c} y_j \operatorname{softmax}(f_\theta(x))_j.$$

Here, $\mathcal{D}$ is a dataset of finite size, with $m = |\mathcal{D}|$. For a specific dataset

$$\mathcal{D}_m = \{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}^c \mid 1 \le i \le m\},$$

we aim to solve the following optimization problem:

$$(\mathcal{P}): \quad \min_{\theta \in \Theta} l(\theta, \mathcal{D}_m) + \frac{\kappa}{2m} \|\theta\|_2^2.$$

Here, $\kappa$ is a regularization weight that penalizes the magnitude of the parameters. The focus of our study is to compare various optimization methods for solving the problem $(\mathcal{P})$.

## 2.2 Algorithms

In this subsection, we outline the optimization methods used to solve the problem $(\mathcal{P})$. These include Gradient Descent (GD), Stochastic Gradient Descent (SGD), and the Adam optimizer. Each method is described below.

### 2.2.1 Gradient Descent

The classic Gradient Descent algorithm iteratively updates the parameters $\theta$ by following the negative gradient of the loss function.

---

**Algorithm 1** Gradient Descent (GD)

---

1: Initialize $\theta_0$, $n \leftarrow 0$, $\eta > 0$, $\kappa \ge 0$ and number of epochs $T$
2: **for** $n = 1, 2, \ldots, T$ **do**
3: $\quad \theta_n = \theta_{n-1}\left(1 - \frac{\eta\kappa}{m}\right) - \frac{\eta}{m}\sum_{i=1}^{m} \nabla_\theta l(\theta_{n-1}, \{x_i, y_i\})$
4: **end for**
5: **return** $\theta_n$

---

However, since the gradient is computed in a single run over the whole dataset, convergence will be pretty slow and the main bottleneck is memory inefficiency. We therefore next present two more advanced algorithms, that are also used in practice, to compare their performance afterward.

### 2.2.2 Stochastic Gradient Descent

Stochastic gradient descent works by studying the Markov chain $\{X_n; n \in \mathbb{N}\}$ defined through its random mapping representation $X_{n+1} = f_{n+1}(X_n)$ where

$$f_{n+1}(\theta) = \theta\left(1 - \frac{\eta\kappa}{m}\right) - \frac{\eta}{b} \sum_{i \in B_{n+1}} \nabla_\theta l(\theta, \{x_i, y_i\}).$$

Here, $b$ is the (fixed) size of the mini-batch $B_{n+1}$, a random subset of $[\![1, m]\!]$, dependent on $n + 1$ as it is reshuffled at each iteration. The functions $f_n$ for $n$ in $\mathbb{N}^*$ are iid copies of a locally Lipschitz random mapping $f \colon \mathbb{R}^d \to \mathbb{R}^d$.[1]

---

[1]Yanlin Qu, Jose Blanchet, and Peter Glynn. *Deep Learning for Computing Convergence Rates of Markov Chains*. 2024. arXiv: 2405.20435 [cs.LG]. URL: https://arxiv.org/abs/2405.20435.

---

**Algorithm 2** Stochastic Gradient Descent (SGD)

---

1: Initialize $X_0$, $b \in \mathbb{N}^*$, $k \leftarrow 0$, $\eta > 0$, $\kappa \geq 0$ and number of epochs $T$
2: **for** $n = 1, 2, \ldots, T$ **do**
3:      Sample random permutation $\pi$ according to $\pi \sim \mathcal{U}(\mathfrak{S}_m)$
4:      Divide $\mathcal{D}_m$ into $N = \lceil m/b \rceil$ mini-batches $B_j^{(n)} = \pi(\llbracket (j-1)b + 1, \min(jb, m) \rrbracket)$ for $j$ in $\llbracket 1, N \rrbracket$
5:      **for** $j = 1, \ldots, N$ **do**
6:          $X_k \leftarrow X_k \left(1 - \dfrac{\eta\kappa}{m}\right) - \dfrac{\eta}{b} \displaystyle\sum_{i \in B_j^{(n)}} \nabla_\theta l(X_k, \{x_i, y_i\})$
7:          $k \leftarrow k + 1$
8:      **end for**
9: **end for**
10: **return** $\theta_k$

---

Note that when one chooses $b = m$, the algorithm becomes deterministic and corresponds exactly to the gradient descent algorithm. In fact, one will only have one mini-batch $\pi(\llbracket 1, m \rrbracket) = \llbracket 1, m \rrbracket$ for all $\pi$ in $\mathfrak{S}_m$. Now, as we've seen in class, momentum is a point that enhances convergence and renders the algorithm less prone to local minima. The following algorithm, Adam, goes beyond the idea of first-order momentum.

### 2.2.3 Adam Optimization

The Adam optimizer extends SGD by incorporating momentum and adaptive learning rates. It tracks the first- and second-order moments of the gradient using moving averages. Bias-corrected estimates $\hat{u}_n$ and $\hat{v}_n$ are computed to ensure stability. The update rule is:

$$\theta_k \leftarrow \theta_k - \eta \, \hat{u}_k \odot \left(\hat{v}_k + \epsilon^2\right)^{-1/2},$$

where $\eta > 0$ is the learning rate, $\epsilon > 0$ is a small constant for numerical stability, and $\kappa \geq 0$ is the regularization weight.

---

**Algorithm 3** Adam

---

1: Initialize $X_0$, $u_0 \leftarrow 0$, $v_0 \leftarrow 0$, $k \leftarrow 0$, $\alpha > 0$, $\beta > 0$, $\epsilon > 0$, $\eta > 0$, $\kappa \geq 0$, $b \in \mathbb{N}^*$, and number of epochs $T$
2: **for** $n = 1, 2, \ldots, T$ **do**
3:      Sample random permutation $\pi$ according to $\pi \sim \mathcal{U}(\mathfrak{S}_m)$
4:      Divide $\mathcal{D}_m$ into $N = \lceil m/b \rceil$ mini-batches $B_j^{(n)} = \pi(\llbracket (j-1)b + 1, \min(jb, m) \rrbracket)$ for $j$ in $\llbracket 1, N \rrbracket$
5:      **for** $j = 1, \ldots, N$ **do**
6:          $g_k \leftarrow \dfrac{1}{b} \displaystyle\sum_{i \in B_j^{(n)}} \nabla_\theta l(X_k, \{x_i, y_i\}) + \dfrac{\kappa}{m} X_k$
7:          $u_k \leftarrow \alpha u_k + (1 - \alpha)g_k$
8:          $v_k \leftarrow \beta v_k + (1 - \beta)g_k \odot g_k$
9:          $\hat{u}_k \leftarrow u_k / \left(1 - \alpha^{k+1}\right)$
10:          $\hat{v}_k \leftarrow v_k / \left(1 - \beta^{k+1}\right)$
11:          $X_k \leftarrow X_k - \eta \, \hat{u}_k \odot \left(\hat{v}_k + \epsilon^2\right)^{-1/2}$
12:          $k \leftarrow k + 1$
13:      **end for**
14: **end for**
15: **return** $\theta_k$

---

# 3   Results

## 3.1   Parameters

The practical values used in this study are summarized in Table 1. Parameters are categorized into Model Parameters, Training Settings, and Optimization Settings for clarity. Mathematical notations are provided where applicable, along with the corresponding values.

| Category | Notation | Value |
|---|---|---|
| **Model Parameters** | | |
| Input Dimension | $p$ | 280 |
| Number of Classes | $c$ | 3 |
| Hidden Layer Dimension | $h$ | 32 |
| Parameter Dimension | $d$ | 9091 |
| Activation Function | $\sigma$ | ReLU |
| **Training Settings** | | |
| Learning Rate | $\eta$ | 0.001 |
| Batch Size | $b$ | 64 |
| Number of Epochs | $T$ | 20 |
| Loss Function | $l(\theta, \mathcal{D})$ | Cross-Entropy Loss |
| Regularization weight | $\kappa$ | 0.7873 |
| Early Stopping Patience | – | 5 epochs |
| Size of Training Set | $m$ | 7873 |
| Train-Test Splitting Factor | – | 0.8 |
| **Optimization Settings** | | |
| Optimizer | – | Adam, SGD, GD (comparative study) |
| Learning Rate Scheduler | StepLR$(\tau, \gamma)$ | StepLR (step size: $\tau = 20$, decay factor: $\gamma = 0.5$) |
| Betas | $(\alpha, \beta)$ | (0.9, 0.999) |
| Numerical Stabilizer | $\epsilon$ | 0.001 |

Table 1: Parameters, mathematical notation, and corresponding values used for training and evaluation. Parameters are grouped into model-specific settings, training configurations, and optimization techniques.

## 3.2   Results

### 3.2.1   Performance Comparison

The performance of the models trained using the SGD and Adam optimizers is summarized in the table below:

| Optimizer | Final Validation Loss | Final Validation Accuracy (%) |
|---|---|---|
| SGD | 1.0303 | 48.96 |
| Adam | 1.0690 | 48.84 |

Table 2: Performance comparison of SGD and Adam optimizers.

### 3.2.2   Convergence Plots

The figures below show the convergence of the training and validation losses and the validation accuracies for both optimizers.
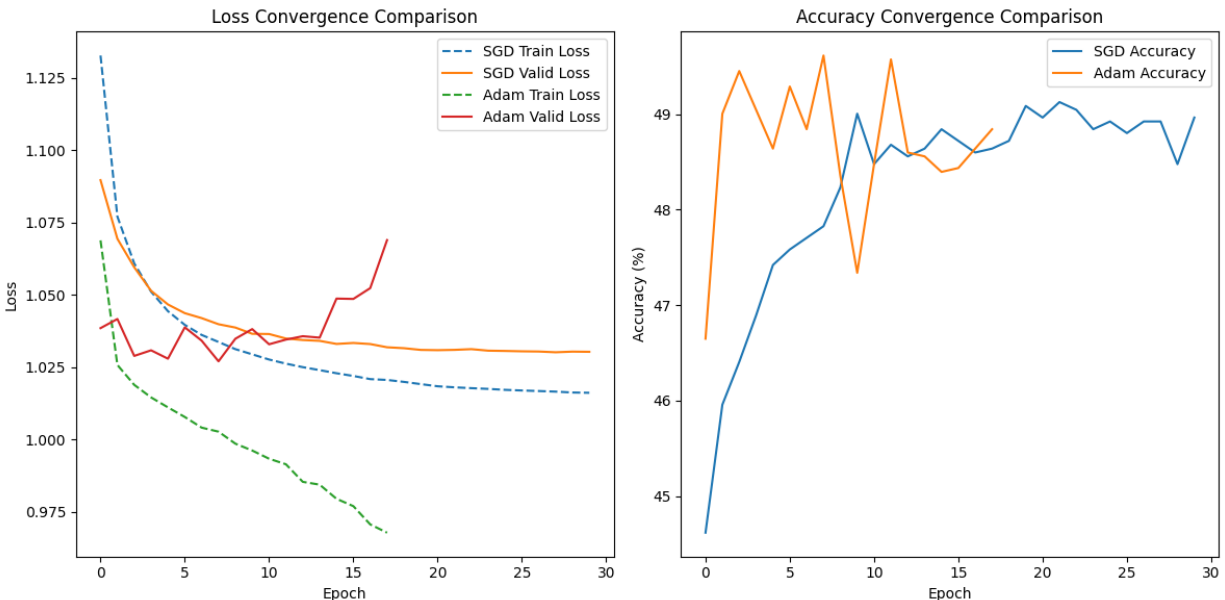
Figure 1: Left: Loss convergence for SGD and Adam optimizers. Right: Accuracy convergence for SGD and Adam optimizers.

# 4    Conclusion

This study evaluated the performance of feedforward neural networks trained with two optimization algorithms, SGD and Adam, to predict football match outcomes. Key results were compared against two benchmark models: (1) the baseline predicting that the HOME team wins, achieving 44% accuracy, and (2) a gradient boosting tree model trained solely on numerical features, yielding 47.5% accuracy on the training set.

The neural network models achieved validation accuracies of 48.96% with SGD and 48.84% with Adam, demonstrating marginal improvement over the benchmarks. Both optimizers displayed similar final performance but revealed notable differences in their behavior:

- SGD achieved slightly higher validation accuracy and lower final validation loss, indicating a more stable convergence to a solution that generalizes better. The simplicity of SGD's fixed learning rate may have contributed to more consistent updates during training, particularly for this problem's dataset and network structure.

- Adam, while competitive, showed a slightly higher final validation loss, which suggests a possible issue with overfitting or sensitivity to hyperparameter settings. Adam's adaptive learning rate mechanism can sometimes cause oscillations or premature convergence in scenarios where the learning rate becomes too aggressive or inconsistent. This could explain its slightly inferior performance compared to SGD in this study.

Despite these challenges, Adam exhibited faster initial convergence, as observed in the loss and accuracy plots. Its ability to adjust learning rates dynamically allowed it to make rapid progress early in training, which is a notable advantage over SGD. This aspect of Adam can be particularly beneficial in cases where computational efficiency or early stopping is critical.

In summary, while SGD proved marginally superior in this specific task due to its stability and generalization, Adam's faster convergence remains advantageous in certain contexts. The results underscore the importance of tailoring optimization strategies to the characteristics of the dataset and model. Future work could involve further hyperparameter tuning, regularization techniques, and hybrid training strategies to fully harness the strengths of both optimizers and achieve better overall performance.

**Appendix: Python3 Code**

```python
1   import pandas as pd
2   import numpy as np
3   from sklearn import model_selection
4   import warnings
5   import matplotlib.pyplot as plt
6   import torch
7   from torch.utils.data import DataLoader, TensorDataset
8   import torch.nn as nn
9   import torch.optim as optim
10
11  warnings.filterwarnings('ignore')
12
13  # Load and preprocess data
14  train_home_team_statistics_df = pd.read_csv('data/Train_Data/train_home_team_statistics_df.
        csv', index_col=0)
15  train_away_team_statistics_df = pd.read_csv('data/Train_Data/train_away_team_statistics_df.
        csv', index_col=0)
16  train_scores = pd.read_csv('data/Y_train_1rknArQ.csv', index_col=0)
17
18  train_home = train_home_team_statistics_df.iloc[:, 2:]
19  train_away = train_away_team_statistics_df.iloc[:, 2:]
20
21  train_home.columns = 'HOME_' + train_home.columns
22  train_away.columns = 'AWAY_' + train_away.columns
23
24  train_data = pd.concat([train_home, train_away], join='inner', axis=1)
25  train_scores = train_scores.loc[train_data.index]
26
27  train_data = train_data.replace({np.inf: np.nan, -np.inf: np.nan})
28  train_scores_indices = train_scores.idxmax(axis=1).map({'HOME_WINS': 0, 'DRAW': 1, '
        AWAY_WINS': 2})
29
30  X_train, X_valid, y_train, y_valid = model_selection.train_test_split(
31      train_data, train_scores_indices, train_size=0.8, random_state=42
32  )
33
34  # Prepare data for PyTorch
35  X_train_tensor = torch.from_numpy(X_train.fillna(0).to_numpy()).float()
36  y_train_tensor = torch.from_numpy(y_train.to_numpy()).long()
37  X_valid_tensor = torch.from_numpy(X_valid.fillna(0).to_numpy()).float()
38  y_valid_tensor = torch.from_numpy(y_valid.to_numpy()).long()
39
40  train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
41  valid_dataset = TensorDataset(X_valid_tensor, y_valid_tensor)
42
43  train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
44  valid_loader = DataLoader(valid_dataset, batch_size=64, shuffle=False)
45
46  device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
47
48  # Define the Neural Network
49  class NeuralNet(nn.Module):
50      def __init__(self, input_size, hidden_size, num_classes):
51          super(NeuralNet, self).__init__()
52          self.fc1 = nn.Linear(input_size, hidden_size)
53          self.relu = nn.ReLU()
54          self.fc2 = nn.Linear(hidden_size, num_classes)
55
56      def forward(self, x):
57          x = self.fc1(x)
58          x = self.relu(x)
59          x = self.fc2(x)
60          return x
61
62  # Define evaluation function
63  def evaluate(model, criterion, loader):
```

```python
64        model.eval()
65        total_loss, correct, total = 0, 0, 0
66        with torch.no_grad():
67            for inputs, labels in loader:
68                inputs, labels = inputs.to(device), labels.to(device)
69                outputs = model(inputs)
70                total_loss += criterion(outputs, labels).item()
71                predicted = outputs.argmax(dim=1)
72                correct += (predicted == labels).sum().item()
73                total += labels.size(0)
74        accuracy = 100 * correct / total
75        return total_loss / len(loader), accuracy
76
77   # Early stopping implementation
78   class EarlyStopping:
79        def __init__(self, patience, verbose=False):
80            self.patience = patience
81            self.verbose = verbose
82            self.best_loss = float('inf')
83            self.counter = 0
84
85        def step(self, current_loss):
86            if current_loss < self.best_loss:
87                self.best_loss = current_loss
88                self.counter = 0
89                return True
90            else:
91                self.counter += 1
92                if self.counter >= self.patience:
93                    if self.verbose:
94                        print("Early stopping triggered.")
95                    return False
96            return True
97
98   # Train the model
99   def train_model(optimizer_name, lr=0.001, weight_decay=1e-5, hidden_size=32, num_epochs=30,
         early_stopping_patience=10):
100       input_size = X_train_tensor.shape[1]
101       num_classes = len(np.unique(y_train))
102
103       model = NeuralNet(input_size, hidden_size, num_classes).to(device)
104       criterion = nn.CrossEntropyLoss()
105
106       if optimizer_name == 'SGD':
107           optimizer = optim.SGD(model.parameters(), lr=lr, weight_decay=weight_decay)
108       elif optimizer_name == 'Adam':
109           optimizer = optim.Adam(model.parameters(), lr=lr, weight_decay=weight_decay)
110       else:
111           raise ValueError("Unknown optimizer")
112
113       scheduler = optim.lr_scheduler.StepLR(optimizer, step_size=20, gamma=0.5)
114       early_stopping = EarlyStopping(patience=early_stopping_patience, verbose=True)
115
116       train_losses, valid_losses, accuracies = [], [], []
117       learning_rates = []
118
119       for epoch in range(num_epochs):
120           model.train()
121           train_loss = 0
122           for inputs, labels in train_loader:
123               inputs, labels = inputs.to(device), labels.to(device)
124
125               outputs = model(inputs)
126               loss = criterion(outputs, labels)
127
128               optimizer.zero_grad()
129               loss.backward()
130               optimizer.step()
```

```
131
132                 train_loss += loss.item()
133
134         train_loss /= len(train_loader)
135         train_losses.append(train_loss)
136
137         valid_loss, valid_accuracy = evaluate(model, criterion, valid_loader)
138         valid_losses.append(valid_loss)
139         accuracies.append(valid_accuracy)
140
141         learning_rates.append(optimizer.param_groups[0]['lr'])
142
143         print(f'Epoch [{epoch+1}/{num_epochs}], Loss: {valid_loss:.4f}, Accuracy: {
                valid_accuracy:.2f}%, Learning Rate: {learning_rates[-1]:.6f}')
144
145         if not early_stopping.step(valid_loss):
146             break
147
148         scheduler.step()
149
150     return train_losses, valid_losses, accuracies, valid_loss, valid_accuracy,
            learning_rates
151
152 # Plot metrics
153 def plot_metrics(train_losses, valid_losses, accuracies, title_suffix=""):
154     plt.figure(figsize=(12, 6))
155     plt.subplot(1, 2, 1)
156     plt.plot(train_losses, label='Train Loss')
157     plt.plot(valid_losses, label='Valid Loss')
158     plt.xlabel('Epoch')
159     plt.ylabel('Loss')
160     plt.title(f'Loss Convergence {title_suffix}')
161     plt.legend()
162
163     plt.subplot(1, 2, 2)
164     plt.plot(accuracies, label='Accuracy')
165     plt.xlabel('Epoch')
166     plt.ylabel('Accuracy (%)')
167     plt.title(f'Accuracy Convergence {title_suffix}')
168     plt.legend()
169     plt.tight_layout()
170     plt.show()
171
172 # Plot learning rate progression
173 def plot_learning_rate(learning_rates, title_suffix=""):
174     plt.figure(figsize=(6, 4))
175     plt.plot(learning_rates, label='Learning Rate')
176     plt.xlabel('Epoch')
177     plt.ylabel('Learning Rate')
178     plt.title(f'Learning Rate Progression {title_suffix}')
179     plt.legend()
180     plt.tight_layout()
181     plt.show()
182
183 # Plot comparison
184 def plot_comparison(sgd_results, adam_results):
185     plt.figure(figsize=(12, 6))
186
187     plt.subplot(1, 2, 1)
188     plt.plot(sgd_results[0], label='SGD Train Loss', linestyle='--')
189     plt.plot(sgd_results[1], label='SGD Valid Loss')
190     plt.plot(adam_results[0], label='Adam Train Loss', linestyle='--')
191     plt.plot(adam_results[1], label='Adam Valid Loss')
192     plt.xlabel('Epoch')
193     plt.ylabel('Loss')
194     plt.title('Loss Convergence Comparison')
195     plt.legend()
196
```

```
197        plt.subplot(1, 2, 2)
198        plt.plot(sgd_results[2], label='SGD Accuracy')
199        plt.plot(adam_results[2], label='Adam Accuracy')
200        plt.xlabel('Epoch')
201        plt.ylabel('Accuracy (%)')
202        plt.title('Accuracy Convergence Comparison')
203        plt.legend()
204
205        plt.tight_layout()
206        plt.show()
207
208    # Train and evaluate models
209    sgd_results = train_model('SGD')
210    adam_results = train_model('Adam')
211
212    # Plot results
213    plot_comparison(sgd_results, adam_results)
214    plot_learning_rate(sgd_results[5], title_suffix="(SGD)")
215    plot_learning_rate(adam_results[5], title_suffix="(Adam)")
216
217    # Final performance
218    print(f"SGD Final Validation Loss: {sgd_results[3]:.4f}, Final Validation Accuracy: {
           sgd_results[4]:.2f}%")
219    print(f"Adam Final Validation Loss: {adam_results[3]:.4f}, Final Validation Accuracy: {
           adam_results[4]:.2f}%")
```

Listing 1: See the repository https://github.com/juliusgraf/challenge_qrt_optim.