

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC VÀ KỸ THUẬT MÁY TÍNH



Mathematics Foundations for Computer Science (CO5097)

ASSIGNMENT

**OPTIMIZATION FOR MACHINE LEARNING
ADAGRAD**

Lecturer: Dr. Nguyễn An Khương
GROUP: 12
Members: 1. Ngô Nhất Toàn - 2570515
2. Trần Hoàng Ân - 2570548
3. Trần Huy Phước - 2570482
4. Lê Thành Đạt - 2370497
5. Phạm Minh Quang - 2570302

Mục lục

1	Nền tảng lý thuyết của tối ưu hóa	3
1.1	Tối ưu hóa và hàm mất mát	3
1.2	Gradient	4
1.3	Ma trận Hessian	4
1.4	Các thách thức chính trong tối ưu hóa	5
1.4.1	Cực tiểu địa phương (Local Minima)	5
1.4.2	Điểm yên ngựa (Saddle Points)	6
1.4.3	Triệt tiêu Gradient (Vanishing Gradient)	6
2	Các Thuật Toán Nền Tảng và Hạn Chế	8
2.1	Giới thiệu bộ dữ liệu RCV1	8
2.1.1	Giới thiệu Tổng quan	8
2.1.2	Biểu diễn Dữ liệu	8
2.1.3	Dữ liệu thưa (Sparse Features) là gì?	9
2.2	Gradient Descent (GD)	9
2.3	Stochastic Gradient Descent (SGD)	12
2.3.1	Stochastic Gradient Descent	12
2.4	Hạn Chế của SGD với Dữ liệu Thưa (Sparse Features)	21
2.4.1	Phân tích Gradient với Dữ liệu Thưa	21
2.4.2	Vấn đề xảy ra với dữ liệu thưa của Learning Rate	21
2.5	Phân tích Dữ liệu Thưa bằng Toán học	22
2.5.1	Gradient trong Binary Classification với BCE Loss	22
2.5.2	Minh họa với ví dụ cụ thể	23
2.5.3	Tần suất Cập nhật và Learning Rate Requirements	23
2.5.4	Sparse Gradient và Ảnh hưởng đến Optimization	24
3	Adagrad (Adaptive Gradient Algorithm) - Giải Pháp	26
3.1	Động lực: Tốc độ học Thích ứng (Motivation)	26
3.2	Tiền điều kiện của Adagrad	26
3.2.1	Giới thiệu	26
3.2.2	Sự biến đổi của hàm lỗi bậc hai	27
3.3	Thuật toán Adagrad	30
3.4	Ví dụ Tính toán Cụ thể cho Adagrad	32

3.4.1	Phần 1: Ví dụ Tính toán "Coordinate-wise" (Từng tọa độ)	32
3.4.2	Phần 2: Lý thuyết Khởi tạo Trọng số (Weight Initialization)	34
3.5	Phân tích Ưu điểm và Nhược điểm	35
3.6	Adagrad trên Python	35
3.6.1	Minh hoạ tối ưu hoá Adagrad trên hàm số	35
3.6.2	Xây dựng thuật toán Adagrad từ công thức	37
3.6.3	Minh hoạ Adagrad trên tập RCV1	39
3.7	So sánh giữa Adagrad và SGD trên tập RCV1	48
3.7.1	Thiết lập Thực nghiệm	48
3.7.2	Kết quả Thực nghiệm	48
3.7.3	So sánh giữa hai thuật toán tối ưu	48
3.8	Bài tập – Chương 11.7 Adagrad (D2L)	49
4	Tổng kết và Hướng phát triển	62
4.1	Tóm tắt	62
4.2	Các hướng cải tiến (Future Work)	62
	Tài liệu tham khảo	63

1 Nền tảng lý thuyết của tối ưu hóa

1.1 Tối ưu hóa và hàm mất mát

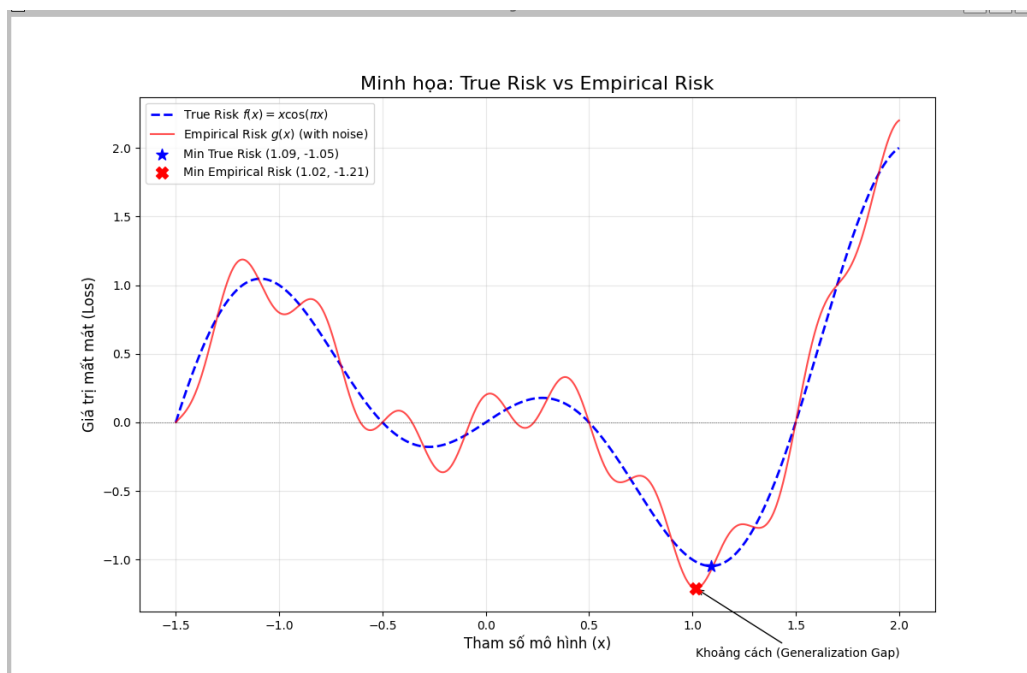
Trong học có giám sát, sai số huấn luyện (hoặc rủi ro thực nghiệm) đo mức độ sai lệch giữa dự đoán của mô hình và nhãn thực tế, được định nghĩa:

$$R_{\text{emp}}[\mathbf{X}, \mathbf{y}, f] = \frac{1}{n} \sum_{i=1}^n l(\mathbf{x}^{(i)}, y^{(i)}, f(\mathbf{x}^{(i)}))$$

Trong đó:

- R_{emp} : Rủi ro thực nghiệm được tính trên tập huấn luyện
- $l(\dots)$: Hàm mất mát đo sự khác biệt giữa dự đoán và nhãn thực tế
- $f(\mathbf{x}^{(i)})$: Dự đoán của mô hình cho đầu vào $\mathbf{x}^{(i)}$
- n : Số lượng mẫu huấn luyện

Mục tiêu tối ưu hóa là cực tiểu hóa R_{emp} (empirical risk), nhưng mục tiêu cuối cùng là tìm mô hình tổng quát hóa tốt với rủi ro thực (true risk), tránh hiện tượng quá khớp khi khoảng cách giữa hai rủi ro quá lớn.



Hình 1: So sánh True Risk và Empirical Risk minh họa khoảng cách tổng quát hóa

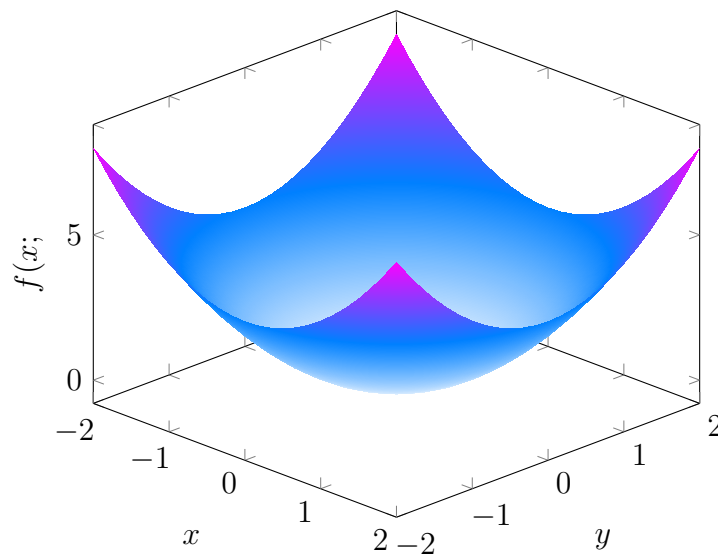
1.2 Gradient

Định nghĩa: Gradient của một hàm đa biến là vector chứa các đạo hàm riêng phần cấp một, chỉ hướng tăng nhanh nhất của hàm số.

Đối với hàm mất mát $L(\mathbf{w})$, gradient được xác định:

$$\nabla L(\mathbf{w}) = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]^T$$

Ví dụ: Với hàm $f(x, y) = x^2 + y^2$, vector gradient là $\nabla f = [2x, 2y]^T$.



Bề mặt paraboloid này có **gradient hướng ra xa tâm**. Trong Gradient Descent, ta di chuyển ngược hướng gradient (xuống dốc) để giảm thiểu hàm mục tiêu. Tại cực tiểu toàn cục $(0, 0)$, gradient bằng 0.

1.3 Ma trận Hessian

Định nghĩa: Ma trận Hessian là ma trận vuông chứa các đạo hàm riêng cấp hai, mô tả độ cong (curvature) của hàm số:

$$H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

Tại điểm tới hạn (nơi $\nabla f = 0$), Hessian phân loại điểm đó:

- H xác định dương ($\det(H) > 0, H_{11} > 0$): Cực tiểu cục bộ
- H xác định âm ($\det(H) < 0, H_{11} < 0$): Cực đại cục bộ
- H không xác định (có cả giá trị riêng dương và âm): Điểm yên ngựa

Ví dụ: Kiểm tra điểm $P(1, 1)$ của hàm $f(x, y) = x^3 + y^3 - 3xy$.

Bước 1 - Đạo hàm bậc nhất:

$$f_x = 3x^2 - 3y, \quad f_y = 3y^2 - 3x$$

Bước 2 - Ma trận Hessian:

$$H(x, y) = \begin{bmatrix} 6x & -3 \\ -3 & 6y \end{bmatrix}$$

Bước 3 - Tại $P(1, 1)$:

$$H(1, 1) = \begin{bmatrix} 6 & -3 \\ -3 & 6 \end{bmatrix}$$

Bước 4 - Phân loại:

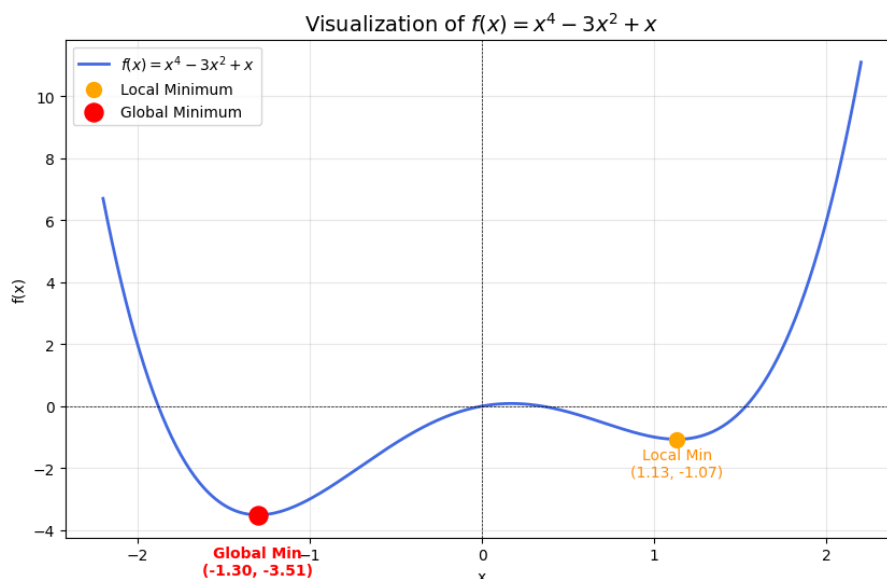
- $D_1 = 6 > 0$
- $D_2 = \det(H) = 36 - 9 = 27 > 0$

Kết luận: Vì $D_1 > 0$ và $D_2 > 0$, ma trận Hessian xác định dương, nên $P(1, 1)$ là điểm cực tiểu cục bộ.

1.4 Các thách thức chính trong tối ưu hóa

1.4.1 Cực tiểu địa phương (Local Minima)

Thuật toán bị mắc kẹt tại điểm có gradient $\nabla L \approx 0$ nhưng chưa đạt cực tiểu toàn cục.



Hình 2: Hàm $f(x) = x^4 - 3x^2 + x$ minh họa cực tiểu địa phương và toàn cục

1.4.2 Điểm yên ngựa (Saddle Points)

Tại điểm yên ngựa, gradient bằng 0 nhưng không phải cực trị—bề mặt cong lên theo một hướng và cong xuống theo hướng khác. Gradient rất nhỏ gần điểm này làm thuật toán di chuyển cực chậm.

Ví dụ phân tích: Hàm $f(x, y) = x^2 - y^2$ tại điểm $(0, 0)$.

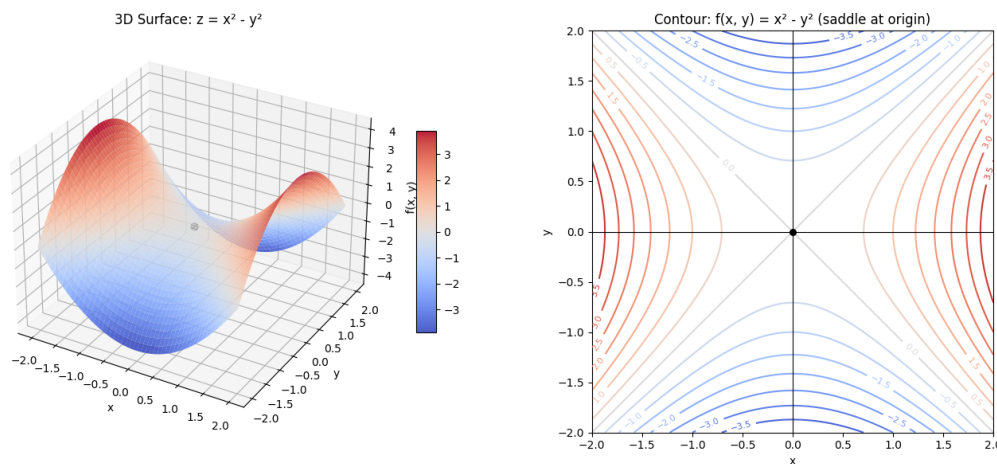
1. **Gradient:** $\nabla f = [2x, -2y]^T$. Tại $(0, 0)$: $\nabla f = \mathbf{0}$ (điểm dừng)

2. **Ma trận Hessian:**

$$H = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

3. **Giá trị riêng:** $\lambda_1 = 2 > 0$, $\lambda_2 = -2 < 0$

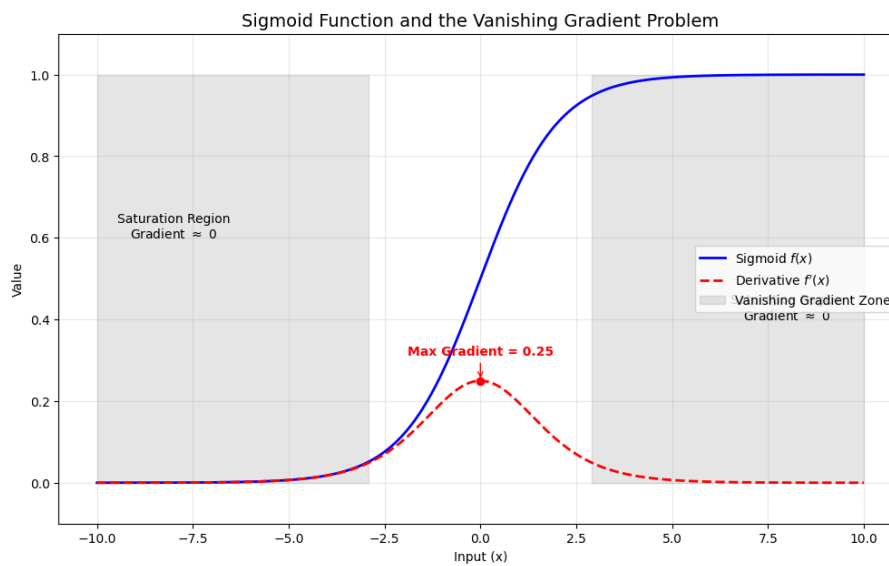
Kết luận: Ma trận không xác định (có cả giá trị riêng dương và âm) xác nhận $(0, 0)$ là điểm yên ngựa.



Hình 3: Hình yên ngựa của $f(x, y) = x^2 - y^2$ với điểm tối hạn tại gốc tọa độ

1.4.3 Triệt tiêu Gradient (Vanishing Gradient)

Gradient giảm dần về 0 khi lan truyền ngược qua nhiều lớp sâu, do tích của nhiều đạo hàm nhỏ (ví dụ: Sigmoid có đạo hàm tối đa 0.25).



Hình 4: Hàm Sigmoid $\sigma(x) = 1/(1 + e^{-x})$ gây triệt tiêu gradient

2 Các Thuật Toán Nền Tảng và Hạn Chế

Phần này trình bày các thuật toán gradient descent nền tảng và những hạn chế khi áp dụng vào các bài toán thực tế, đặc biệt trong bối cảnh dữ liệu thưa (sparse data) như đối với bộ dữ liệu RCV1.

2.1 Giới thiệu bộ dữ liệu RCV1

2.1.1 Giới thiệu Tổng quan

RCV1 (Reuters Corpus Volume 1) là bộ dữ liệu chuẩn (benchmark dataset) được sử dụng rộng rãi trong nghiên cứu Machine Learning, đặc biệt cho các bài toán phân loại văn bản và phân tích dữ liệu thưa (sparse data analysis).

Bộ dữ liệu này được thu thập và công bố bởi Reuters, bao gồm hơn 800,000 bài báo. tiếng Anh được xuất bản bởi hãng thông tấn Reuters trong khoảng thời gian từ tháng 8/1996 đến 8/1997. Mỗi bài báo được gán nhãn với một hoặc nhiều chủ đề từ hệ thống phân loại dạng phân cấp của Reuters.

Đặc điểm	Training Set	Test Set	Tổng
Số documents	23,149	781,265	804,414
Số features (terms)	47,236	47,236	47,236
Số topics/categories	103	103	103
Kích thước ma trận	$23,149 \times 47,236$	$781,265 \times 47,236$	-
Số phần tử khác 0	$\sim 1,500,000$	$\sim 49,000,000$	$\sim 50,500,000$
Sparsity (độ thưa)	99.86%	99.87%	$\sim 99.87\%$

Bảng 1: Thống kê đặc trưng của bộ dữ liệu RCV1

2.1.2 Biểu diễn Dữ liệu

Feature Representation: Bag-of-Words (BoW) với **TF-IDF weighting**.

Mỗi document d được biểu diễn bởi vector:

$$\mathbf{x}_d = [x_1, x_2, \dots, x_{47236}]$$

trong đó:

$$x_i = \text{TF-IDF}(\text{term}_i, d) = \text{TF}(\text{term}_i, d) \times \text{IDF}(\text{term}_i)$$

với:

- $TF(term_i, d)$: Term Frequency — tần suất xuất hiện của từ i trong document d
- $IDF(term_i) = \log \frac{N}{df(term_i)}$: Inverse Document Frequency
 - N : Tổng số documents
 - $df(term_i)$: Số documents chứa $term_i$

Đặc điểm sparse:

- Mỗi document thường chỉ chứa khoảng **200–400 từ** khác nhau.
- Trong 47,236 features, chỉ khoảng **0.5–1%** có giá trị khác 0 cho mỗi document.
- Phần lớn (khoảng **99%**) các phần tử trong vector có giá trị bằng 0.

2.1.3 Dữ liệu thưa (Sparse Features) là gì?

Dữ liệu thưa là dữ liệu mà hầu hết các giá trị là 0, hay nói cách khác, các **features chỉ xuất hiện không thường xuyên** (infrequent features). Đây là đặc trưng phổ biến trong nhiều bài toán thực tế:

- **Xử lý ngôn ngữ tự nhiên (NLP)**: Trong mô hình Bag-of-Words, từ “preconditioning” xuất hiện ít hơn rất nhiều so với từ “learning”.
- **Quảng cáo tính toán (Computational Advertising)**: Các sản phẩm/quảng cáo cụ thể chỉ được một nhóm nhỏ người dùng quan tâm.
- **Hệ thống gợi ý (Collaborative Filtering)**: Mỗi user chỉ tương tác với một phần rất nhỏ trong tổng số items.

2.2 Gradient Descent (GD)

Ý tưởng: Cập nhật tham số mô hình dựa trên gradient của hàm mất mát tính trên toàn bộ tập dữ liệu training.

Công thức cập nhật:

$$\mathbf{x} \leftarrow \mathbf{x} - \eta \nabla f(\mathbf{x})$$

Trong đó:

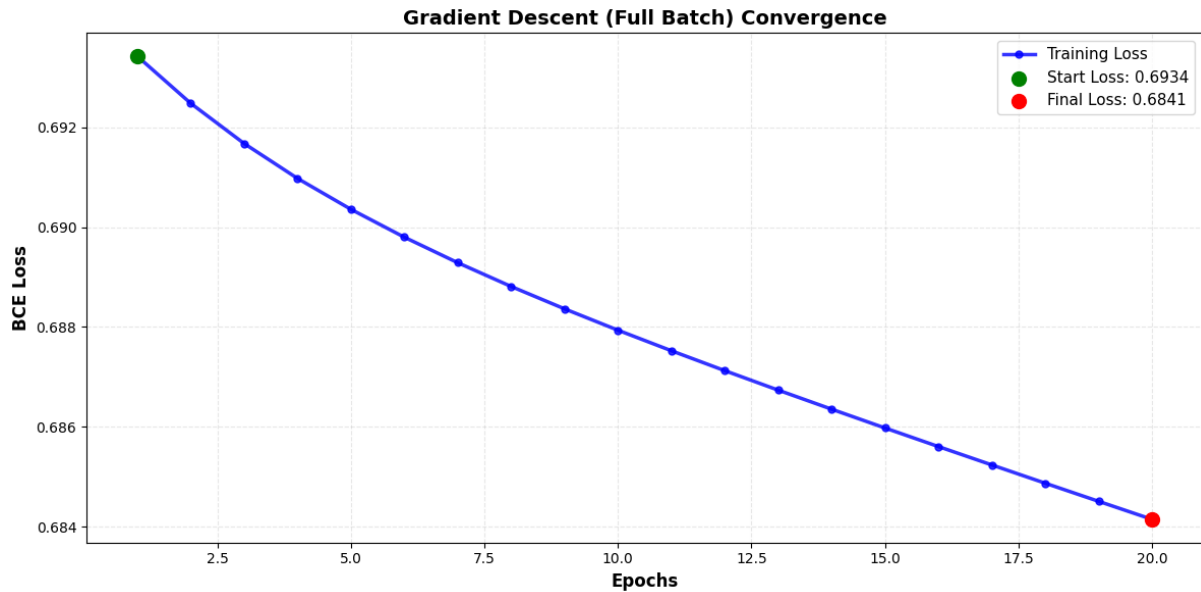
- \mathbf{x} là vector tham số cần tối ưu
- η là learning rate (tốc độ học)
- $\nabla f(\mathbf{x})$ là gradient của hàm mục tiêu tại điểm \mathbf{x}

Vấn đề:

- Chi phí tính toán cực lớn: phải duyệt qua toàn bộ n mẫu để tính gradient cho mỗi bước cập nhật duy nhất.
- Tốn bộ nhớ khi n lớn (như RCV1 với $804,414$ mẫu \times $47,236$ features).
- Tốc độ cập nhật chậm (1 lần/epoch), không tận dụng được random fluctuation.

```
1 def train_gd_visualize(X_sparse, y_tensor, lr=0.5, epochs=20):
2     W = torch.normal(0, 0.01, size=(n_features, 1),
3         requires_grad=True, device=device)
4     b = torch.zeros(1, requires_grad=True, device=device)
5     criterion = nn.BCEWithLogitsLoss()
6
7     epoch_losses = []
8
9     print(f"\nStarting Training Gradient Descent (Full Batch)...")
10    print(f"Total samples: {n_samples} | Updates per epoch: 1")
11    print("="*70)
12
13    start_time = time.time()
14
15    for epoch in range(epochs):
16        epoch_start = time.time()
17
18        linear_out = torch.sparse.mm(X_sparse, W) + b
19
20        loss = criterion(linear_out, y_tensor)
21
22        loss.backward()
```

```
23     # Update parameters (GD update)
24     with torch.no_grad():
25         W -= lr * W.grad
26         b -= lr * b.grad
27         W.grad.zero_()
28         b.grad.zero_()
29
30     current_loss = loss.item()
31     epoch_losses.append(current_loss)
32
33     epoch_time = time.time() - epoch_start
34     if epoch % 5 == 0 or epoch == epochs - 1:
35         print(f"Epoch {epoch+1:3d}/{epochs} | Loss:
36               {current_loss:.6f} | "
37               f"Time: {epoch_time:.2f}s")
38
39     if current_loss < 1e-5:
40         print(f"\n Converged at epoch {epoch+1}!")
41         break
42
43     total_time = time.time() - start_time
44     print("="*70)
45     print(f" Training Complete!")
46     print(f" Total time: {total_time:.2f}s")
47     print(f" Avg time/epoch: {total_time/len(epoch_losses):.2f}s")
48     print(f" Total updates: {len(epoch_losses)}")
49     print(f" Final loss: {epoch_losses[-1]:.6f}")
50
51     return epoch_losses
```



Hình 5: Gradient Descent loss sau 20 epoches

2.3 Stochastic Gradient Descent (SGD)

Ý tưởng chung: Thay vì tính gradient trên toàn bộ dữ liệu như GD, SGD và các biến thể của nó cập nhật tham số dựa trên gradient của một phần nhỏ dữ liệu được chọn ngẫu nhiên. Điều này giúp tăng tốc độ tính toán và khả năng tổng quát hóa.

2.3.1 Stochastic Gradient Descent

Công thức cập nhật:

$$\mathbf{x} \leftarrow \mathbf{x} - \eta \nabla f_i(\mathbf{x}) \quad (1)$$

Trong đó:

- \mathbf{x} là vector tham số cần tối ưu
- $\eta > 0$ là learning rate (tốc độ học)
- $\nabla f_i(\mathbf{x})$ là gradient của hàm loss tại mẫu thứ i
- $i \in \{1, \dots, n\}$ được chọn ngẫu nhiên từ tập training

Đặc điểm:

- Cập nhật rất nhanh (chỉ xử lý 1 mẫu/bước).
- Gradient dao động mạnh (high variance) do chỉ dựa vào 1 mẫu.

- Đường đi đến minimum có dạng “zigzag”.
 - Khả năng thoát local minima nhờ “nhiều” ngẫu nhiên.
-

```
1 import torch
2 import time
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from sklearn.datasets import fetch_rcv1
6 import torch.nn as nn
7 from sklearn.datasets import clear_data_home
8 clear_data_home()
9
10 # 1. Setup & Load Data
11 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
12 print(f"Device: {device}")
13
14 print("Downloading/Loading RCV1...")
15 rcv1 = fetch_rcv1(subset='train')
16 n_samples, n_features = rcv1.data.shape
17
18 # Xu lý nhãn
19 topic_counts = rcv1.target.sum(axis=0)
20 most_frequent_topic_idx = np.argmax(topic_counts)
21 y_full = rcv1.target[:, most_frequent_topic_idx].toarray().ravel()
22 y_cpu = torch.tensor(y_full, dtype=torch.float32).view(-1, 1)
23
24 # 2. Hàm Train SGD có ghi Log và In tiến độ
25 def train_sgd_visualize(X_scipy, y_cpu, lr=0.01, epochs=1):
26     # Khởi tạo tham số
27     W = torch.normal(0, 0.01, size=(n_features, 1),
28                       requires_grad=True, device=device)
29     b = torch.zeros(1, requires_grad=True, device=device)
30     criterion = nn.BCEWithLogitsLoss()
31
32     step_losses = []
```

32

```
33     print(f"Bắt đầu Training Pure SGD (Total samples:  
        {n_samples})...")
```

```
34     start_time = time.time()
```

35

```
36     for epoch in range(epochs):
```

```
37         indices = np.random.permutation(n_samples)
```

38

```
39         for i, idx in enumerate(indices):
```

```
40             # Lấy 1 mẫu (Batch size = 1)
```

```
41             X_sample = torch.tensor(X_scipy[idx].toarray(),  
                                     dtype=torch.float32).to(device)
```

```
42             y_sample = y_cpu[idx].to(device).view(1, 1)
```

43

```
44             # Forward & Loss
```

```
45             linear_out = X_sample @ W + b
```

```
46             loss = criterion(linear_out, y_sample)
```

47

```
48             # Backward & Update
```

```
49             loss.backward()
```

```
50             with torch.no_grad():
```

```
51                 W -= lr * W.grad
```

```
52                 b -= lr * b.grad
```

```
53                 W.grad.zero_()
```

```
54                 b.grad.zero_()
```

55

```
56             # Lưu lại loss
```

```
57             current_loss = loss.item()
```

```
58             step_losses.append(current_loss)
```

59

```
60             # --- PHẦN BẠN CẦN: IN TIẾN ĐỘ& LOSS ---
```

```
61             # In ra mỗi 1000 samples để theo dõi
```

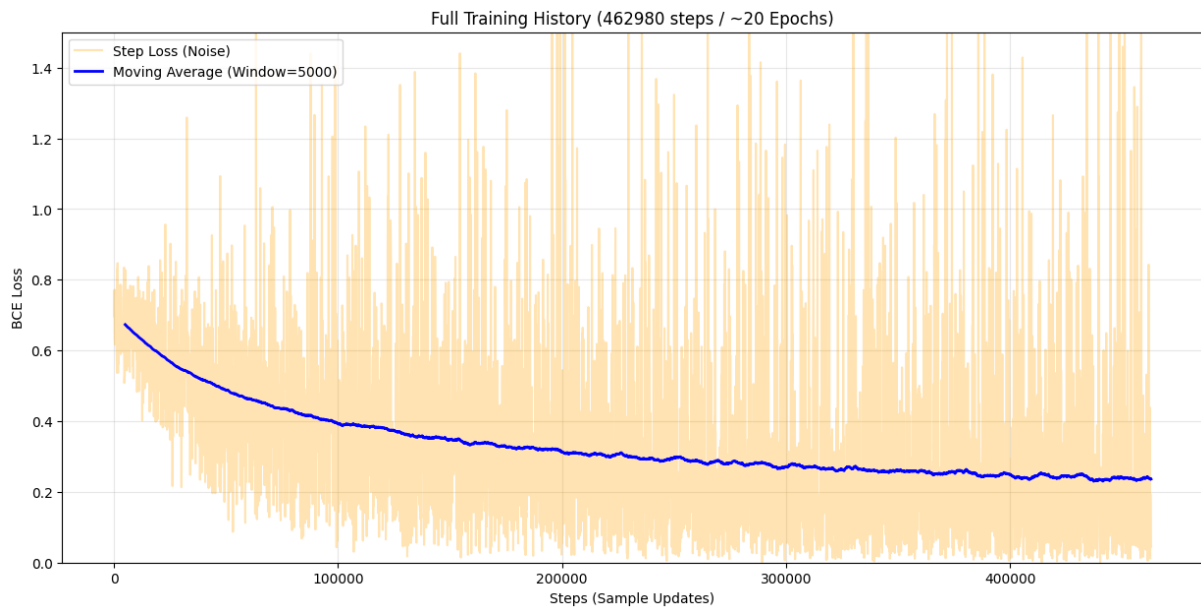
```
62             if i % 5000 == 0:
```

```
63                 print(f"Epoch {epoch+1} | Sample {i}/{n_samples} |  
                    Current Loss: {current_loss:.4f}")
```

```
64
65     print(f"--> Kết thúc Epoch {epoch+1}")
66
67     print(f"Tổng thời gian: {time.time() - start_time:.2f}s")
68     return step_losses
69
70 # 3. Hàm Visualize (Vẽ biểu đồ Nhiễu vs Xu hướng)
71 def plot_full_training_history(losses, window_size=1000,
72                               downsample_rate=100):
73     """
74     Vẽ toàn bộ quá trình huấn luyện bằng cách giảm số lượng điểm vẽ
75     (Downsampling).
76     Input:
77         losses: List chứa toàn bộ history (khoảng 460k điểm)
78         window_size: Kích thước của số Moving Average
79         downsample_rate: Vẽ 1 điểm cho mỗi n bước (đề biểu đồ nhẹ hơn)
80     """
81     plt.figure(figsize=(15, 7))
82
83     total_steps = len(losses)
84
85     # 1. Vẽ Loss thực tế (Downsampling để không bị dày đặc quá)
86     # Lấy bước nhảy là downsample_rate
87     steps = range(0, total_steps, downsample_rate)
88     sampled_losses = [losses[i] for i in steps]
89
90     plt.plot(steps, sampled_losses, color='orange', alpha=0.3,
91             label='Step Loss (Noise)')
92
93     # 2. Vẽ Moving Average cho TOÀN BỘ quá trình
94     # Tính trên dữ liệu gốc để chính xác, sau đó mới vẽ
95     if total_steps > window_size:
96         # Tính Moving Average
97         moving_avg = np.convolve(losses,
98                                 np.ones(window_size)/window_size, mode='valid')
```



```
95
96     # Trục x cho moving average (căn chỉnh về cuối của số)
97     ma_x_axis = range(window_size-1, total_steps)
98
99     # Downsample đường Moving Average để vẽ cho nhanh (tùy chọn)
100    # Vì đường này mượt nên có thể vẽ hết hoặc downsample ít hơn
101    plt.plot(ma_x_axis[::downsample_rate],
102             moving_avg[::downsample_rate],
103             color='blue', linewidth=2, label=f'Moving Average
104             (Window={window_size})')
105
106    plt.title(f'Full Training History ({total_steps} steps /
107            ~{total_steps//23149} Epochs)')
108    plt.xlabel('Steps (Sample Updates)')
109    plt.ylabel('BCE Loss')
110    plt.legend()
111    plt.grid(True, alpha=0.3)
112
113    # Giới hạn trục Y để loại bỏ các điểm nhiễu quá lớn (outliers)
114    plt.ylim(0, 1.5)
115
116    plt.show()
117
118    # --- SỬ DỤNG ---
119    print("Đang vẽ biểu đồ toàn bộ quá trình...")
120    # window_size lớn hơn (5000) để làm mượt đường xu hướng trên tập dữ
121    # liệu lớn
122    loss_history = train_sgd_visualize(rcv1.data, y_cpu, lr=0.01,
123                                     epochs=20)
124
125    plot_full_training_history(loss_history, window_size=5000,
126                              downsample_rate=100)
127
128    # --- CHẠY TRAINING VÀ VẼ ---
129    # Mình để LR nhỏ (0.01) để loss đỡ bị nhảy quá cao
```



Hình 6: Loss của SGD sau 20 epoches

```

1 Device: cuda
2 Downloading/Loading RCV1...
3 Đang vẽ biểu đồtoàn bộ quá trình...
4 Bắt đầu Training Pure SGD (Total samples: 23149)...
5 Epoch 1 | Sample 0/23149 | Current Loss: 0.6814
6 Epoch 1 | Sample 5000/23149 | Current Loss: 0.5973
7 Epoch 1 | Sample 10000/23149 | Current Loss: 0.5429
8 Epoch 1 | Sample 15000/23149 | Current Loss: 0.5499
9 Epoch 1 | Sample 20000/23149 | Current Loss: 0.7225
10 --> Kết thúc Epoch 1
11 Epoch 2 | Sample 0/23149 | Current Loss: 0.5398
12 Epoch 2 | Sample 5000/23149 | Current Loss: 0.5941
13 Epoch 2 | Sample 10000/23149 | Current Loss: 0.4737
14 Epoch 2 | Sample 15000/23149 | Current Loss: 0.3848
15 Epoch 2 | Sample 20000/23149 | Current Loss: 0.4982
16 --> Kết thúc Epoch 2
17 Epoch 3 | Sample 0/23149 | Current Loss: 0.9351
18 Epoch 3 | Sample 5000/23149 | Current Loss: 0.2774
19 Epoch 3 | Sample 10000/23149 | Current Loss: 0.3293
20 Epoch 3 | Sample 15000/23149 | Current Loss: 0.5129
21 Epoch 3 | Sample 20000/23149 | Current Loss: 0.3087

```

```
22 --> Kết thúc Epoch 3
23 Epoch 4 | Sample 0/23149 | Current Loss: 0.3378
24 Epoch 4 | Sample 5000/23149 | Current Loss: 0.5450
25 Epoch 4 | Sample 10000/23149 | Current Loss: 0.4627
26 Epoch 4 | Sample 15000/23149 | Current Loss: 0.7220
27 Epoch 4 | Sample 20000/23149 | Current Loss: 0.5100
28 --> Kết thúc Epoch 4
29 Epoch 5 | Sample 0/23149 | Current Loss: 0.2921
30 Epoch 5 | Sample 5000/23149 | Current Loss: 0.6749
31 Epoch 5 | Sample 10000/23149 | Current Loss: 0.3351
32 Epoch 5 | Sample 15000/23149 | Current Loss: 0.3625
33 Epoch 5 | Sample 20000/23149 | Current Loss: 0.2539
34 --> Kết thúc Epoch 5
35 Epoch 6 | Sample 0/23149 | Current Loss: 0.0771
36 Epoch 6 | Sample 5000/23149 | Current Loss: 0.1250
37 Epoch 6 | Sample 10000/23149 | Current Loss: 0.2944
38 Epoch 6 | Sample 15000/23149 | Current Loss: 0.5138
39 Epoch 6 | Sample 20000/23149 | Current Loss: 0.5287
40 --> Kết thúc Epoch 6
41 Epoch 7 | Sample 0/23149 | Current Loss: 0.4102
42 Epoch 7 | Sample 5000/23149 | Current Loss: 0.2454
43 Epoch 7 | Sample 10000/23149 | Current Loss: 0.1044
44 Epoch 7 | Sample 15000/23149 | Current Loss: 0.0428
45 Epoch 7 | Sample 20000/23149 | Current Loss: 0.4702
46 --> Kết thúc Epoch 7
47 Epoch 8 | Sample 0/23149 | Current Loss: 0.2686
48 Epoch 8 | Sample 5000/23149 | Current Loss: 0.1768
49 Epoch 8 | Sample 10000/23149 | Current Loss: 0.1311
50 Epoch 8 | Sample 15000/23149 | Current Loss: 0.2627
51 Epoch 8 | Sample 20000/23149 | Current Loss: 0.9129
52 --> Kết thúc Epoch 8
53 Epoch 9 | Sample 0/23149 | Current Loss: 0.4589
54 Epoch 9 | Sample 5000/23149 | Current Loss: 0.1154
55 Epoch 9 | Sample 10000/23149 | Current Loss: 0.4845
56 Epoch 9 | Sample 15000/23149 | Current Loss: 0.1469
```

```
57 Epoch 9 | Sample 20000/23149 | Current Loss: 0.2182
58 --> Kết thúc Epoch 9
59 Epoch 10 | Sample 0/23149 | Current Loss: 0.2021
60 Epoch 10 | Sample 5000/23149 | Current Loss: 0.3154
61 Epoch 10 | Sample 10000/23149 | Current Loss: 0.3660
62 Epoch 10 | Sample 15000/23149 | Current Loss: 0.3271
63 Epoch 10 | Sample 20000/23149 | Current Loss: 0.4100
64 --> Kết thúc Epoch 10
65 Epoch 11 | Sample 0/23149 | Current Loss: 0.1486
66 Epoch 11 | Sample 5000/23149 | Current Loss: 0.3888
67 Epoch 11 | Sample 10000/23149 | Current Loss: 0.1130
68 Epoch 11 | Sample 15000/23149 | Current Loss: 0.4094
69 Epoch 11 | Sample 20000/23149 | Current Loss: 0.4013
70 --> Kết thúc Epoch 11
71 Epoch 12 | Sample 0/23149 | Current Loss: 0.0717
72 Epoch 12 | Sample 5000/23149 | Current Loss: 0.1492
73 Epoch 12 | Sample 10000/23149 | Current Loss: 0.2155
74 Epoch 12 | Sample 15000/23149 | Current Loss: 0.4325
75 Epoch 12 | Sample 20000/23149 | Current Loss: 1.0979
76 --> Kết thúc Epoch 12
77 Epoch 13 | Sample 0/23149 | Current Loss: 0.0481
78 Epoch 13 | Sample 5000/23149 | Current Loss: 0.1407
79 Epoch 13 | Sample 10000/23149 | Current Loss: 0.1196
80 Epoch 13 | Sample 15000/23149 | Current Loss: 0.2122
81 Epoch 13 | Sample 20000/23149 | Current Loss: 0.3492
82 --> Kết thúc Epoch 13
83 Epoch 14 | Sample 0/23149 | Current Loss: 0.2175
84 Epoch 14 | Sample 5000/23149 | Current Loss: 0.0829
85 Epoch 14 | Sample 10000/23149 | Current Loss: 0.1098
86 Epoch 14 | Sample 15000/23149 | Current Loss: 0.1257
87 Epoch 14 | Sample 20000/23149 | Current Loss: 0.0178
88 --> Kết thúc Epoch 14
89 Epoch 15 | Sample 0/23149 | Current Loss: 0.5082
90 Epoch 15 | Sample 5000/23149 | Current Loss: 0.2699
91 Epoch 15 | Sample 10000/23149 | Current Loss: 0.2530
```



```
92 Epoch 15 | Sample 15000/23149 | Current Loss: 0.1322
93 Epoch 15 | Sample 20000/23149 | Current Loss: 0.1598
94 --> Kết thúc Epoch 15
95 Epoch 16 | Sample 0/23149 | Current Loss: 0.3083
96 Epoch 16 | Sample 5000/23149 | Current Loss: 0.7078
97 Epoch 16 | Sample 10000/23149 | Current Loss: 0.0097
98 Epoch 16 | Sample 15000/23149 | Current Loss: 0.0638
99 Epoch 16 | Sample 20000/23149 | Current Loss: 0.1450
100 --> Kết thúc Epoch 16
101 Epoch 17 | Sample 0/23149 | Current Loss: 0.0638
102 Epoch 17 | Sample 5000/23149 | Current Loss: 0.0502
103 Epoch 17 | Sample 10000/23149 | Current Loss: 0.1105
104 Epoch 17 | Sample 15000/23149 | Current Loss: 0.0350
105 Epoch 17 | Sample 20000/23149 | Current Loss: 0.1212
106 --> Kết thúc Epoch 17
107 Epoch 18 | Sample 0/23149 | Current Loss: 0.2561
108 Epoch 18 | Sample 5000/23149 | Current Loss: 0.0412
109 Epoch 18 | Sample 10000/23149 | Current Loss: 0.2268
110 Epoch 18 | Sample 15000/23149 | Current Loss: 0.0755
111 Epoch 18 | Sample 20000/23149 | Current Loss: 0.9352
112 --> Kết thúc Epoch 18
113 Epoch 19 | Sample 0/23149 | Current Loss: 0.0202
114 Epoch 19 | Sample 5000/23149 | Current Loss: 0.0389
115 Epoch 19 | Sample 10000/23149 | Current Loss: 0.2170
116 Epoch 19 | Sample 15000/23149 | Current Loss: 0.0901
117 Epoch 19 | Sample 20000/23149 | Current Loss: 0.1243
118 --> Kết thúc Epoch 19
119 Epoch 20 | Sample 0/23149 | Current Loss: 0.2846
120 Epoch 20 | Sample 5000/23149 | Current Loss: 0.1579
121 Epoch 20 | Sample 10000/23149 | Current Loss: 0.0280
122 Epoch 20 | Sample 15000/23149 | Current Loss: 0.0903
123 Epoch 20 | Sample 20000/23149 | Current Loss: 0.0665
124 --> Kết thúc Epoch 20
125 Tổng thời gian: 496.02s
```

2.4 Hạn Chế của SGD với Dữ liệu Thưa (Sparse Features)

Trong bộ dữ liệu RCV1 đang sử dụng:

- Input vector \mathbf{x} có kích thước 47,236 chiều (số từ vựng).
- Tuy nhiên, mỗi văn bản chỉ chứa vài trăm từ khác 0 \rightarrow độ thưa (sparsity) rất cao.
- Một số từ phổ biến (common features) xuất hiện trong hầu hết văn bản.
- Nhiều từ hiếm (rare/infrequent features) chỉ xuất hiện trong một vài văn bản.

2.4.1 Phân tích Gradient với Dữ liệu Thưa

Xét mô hình tuyến tính đơn giản:

$$\hat{y} = \mathbf{w}^\top \mathbf{x}$$

Khi tính gradient của hàm mất mát L theo trọng số w_i :

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \cdot x_i$$

Hệ quả quan trọng:

Loại Feature	Ví dụ	Tần suất xuất hiện	Tham số
Phổ biến	the, is, a	Hầu hết documents	w_{common}
Hiếm	serendipity, quixotic	Rất ít documents	w_{rare}

Loại feature	Tần suất xuất hiện	Gradient	Tình trạng cập nhật
Phổ biến (frequent)	$x_i \neq 0$ thường xuyên	Khác 0 liên tục	w_i được cập nhật liên tục
Hiếm (infrequent)	$x_i = 0$ hầu hết các mẫu	Bằng 0 phần lớn thời gian	w_i hiếm khi được cập nhật

2.4.2 Vấn đề xảy ra với dữ liệu thưa của Learning Rate

Vấn đề xảy ra với dữ liệu thưa:

1. Với features phổ biến (w_{common}):

- Nhận gradient liên tục qua mọi iteration.
- Learning rate giảm dần \rightarrow hội tụ ổn định đến giá trị tối ưu.
- Hoạt động ổn định.

2. Với features hiếm (w_{rare}):

- Chỉ nhận gradient khác 0 khi feature xuất hiện (rất hiếm).
- Khi feature cuối cùng xuất hiện, learning rate đã giảm quá nhiều.
- Không còn đủ “lực” để cập nhật w_i một cách có ý nghĩa.
- Chưa hội tụ dù training đã kết thúc.

Tóm lại: “Tốc độ học giảm quá chậm đối với các đặc trưng xuất hiện thường xuyên, hoặc quá nhanh đối với các đặc trưng xuất hiện không thường xuyên.”

2.5 Phân tích Dữ liệu Thừa bằng Toán học

Phần này làm rõ cơ chế toán học khiến SGD với learning rate cố định không hiệu quả trên dữ liệu thừa, dẫn đến nhu cầu phát triển các thuật toán adaptive.

2.5.1 Gradient trong Binary Classification với BCE Loss

Xét bài toán binary classification với logistic regression:

$$\hat{y} = \sigma(\mathbf{w}^\top \mathbf{x}) = \sigma\left(\sum_{j=1}^d w_j x_j\right)$$

trong đó $\sigma(z) = \frac{1}{1+e^{-z}}$ là sigmoid activation function.

Binary Cross-Entropy Loss cho một mẫu (\mathbf{x}, y) với $y \in \{0, 1\}$:

$$\ell(\mathbf{w}) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

Gradient theo từng tham số w_j :

Áp dụng chain rule:

$$\frac{\partial \ell}{\partial w_j} = \frac{\partial \ell}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w_j}$$

với $z = \mathbf{w}^\top \mathbf{x}$.

Sau khi tính toán, ta có:

$$\frac{\partial \ell}{\partial w_j} = (\hat{y} - y) \cdot x_j$$

Trong đó:

- $(\hat{y} - y)$: Prediction error, phụ thuộc vào chất lượng dự đoán của mô hình.
- $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{x}) \in (0, 1)$: Xác suất dự đoán.
- $y \in \{0, 1\}$: Nhãn thực tế.
- x_j : Giá trị feature thứ j trong mẫu hiện tại.

Quan sát then chốt:

$$x_j = 0 \Rightarrow \frac{\partial \ell}{\partial w_j} = 0$$

\Rightarrow Tham số w_j **không nhận gradient** từ mẫu này, do đó **không được cập nhật**.

Lưu ý quan trọng: Kết quả này **độc lập** với loss function. Miễn là model tuyến tính $\hat{y} = f(\mathbf{w}^\top \mathbf{x})$, gradient luôn tỷ lệ với x_j .

2.5.2 Minh họa với ví dụ cụ thể

Ví dụ: Xét mini-batch chứa văn bản về doanh nghiệp (label $y = 1$ cho CCAT).

Giả sử sau tokenization và TF-IDF encoding:

Feature	TF-IDF	x_j	Gradient	Cập nhật SGD
“company”	cao	$x_{\text{company}} = 0.8$	$\neq 0$	Cập nhật
“merger”	trung bình	$x_{\text{merger}} = 0.3$	$\neq 0$	Cập nhật
“reuters”	0 (không có)	$x_{\text{reuters}} = 0$	0	Không cập nhật
“eigenvalue”	0 (không có)	$x_{\text{eigenvalue}} = 0$	0	Không cập nhật

Hệ quả:

- **Features phổ biến** (“company”, “merger”): nhận gradient thường xuyên \rightarrow học nhanh.
- **Features hiếm** (“eigenvalue”): gradient = 0 ở hầu hết bước cập nhật \rightarrow học rất chậm.

2.5.3 Tần suất Cập nhật và Learning Rate Requirements

Giả sử trong RCV1 training set (23,149 documents):

- Feature “said”: xuất hiện trong $\sim 15,000$ documents (65%).
- Feature “preconditioning”: xuất hiện trong ~ 5 documents (0.02%).

Sau 1,000 iterations SGD (random sampling):

Feature	Số lần nhận gradient $\neq 0$	Magnitude tích lũy	LR lý tưởng
“said”	~ 650	Rất lớn	Nhỏ (~ 0.001)
“preconditioning”	~ 0 (có thể = 0)	Gần 0	Lớn (~ 0.1)

Dilemma với learning rate cố định η :

$$\eta_{\text{SGD}} = \begin{cases} \text{Nhỏ (0.001)} & \Rightarrow w_{\text{common}} \text{ hội tụ ổn định, nhưng } w_{\text{rare}} \text{ học quá chậm hoặc không học} \\ \text{Lớn (0.1)} & \Rightarrow w_{\text{rare}} \text{ có cơ hội học, nhưng } w_{\text{common}} \text{ dao động mạnh / diverge} \end{cases}$$

Minh họa cập nhật SGD:

$$w_j^{(t+1)} = w_j^{(t)} - \eta \cdot (\hat{y}^{(t)} - y^{(t)}) \cdot x_j^{(t)}$$

- Nếu $x_j^{(t)} = 0$ (feature không xuất hiện): $w_j^{(t+1)} = w_j^{(t)}$ (không thay đổi).
- Feature hiếm: $x_j \neq 0$ chỉ trong $\sim 0.02\%$ iterations \rightarrow tổng cộng cập nhật $\ll 1,000$ lần.
- Feature phổ biến: $x_j \neq 0$ trong $\sim 65\%$ iterations \rightarrow cập nhật ~ 650 lần.

\Rightarrow Chênh lệch 30,000x về số lần cập nhật!

2.5.4 Sparse Gradient và Ảnh hưởng đến Optimization

Định nghĩa Sparse Gradient:

Với một document trong RCV1 (trung bình 67 features khác 0 / 47,236 features):

$$\nabla \ell(\mathbf{w}) = \begin{bmatrix} g_1 \\ 0 \\ 0 \\ g_4 \\ \vdots \\ 0 \end{bmatrix} \quad \text{với} \quad \|\nabla \ell\|_0 \approx 67 \ll 47236$$

Trong đó:

- $\|\nabla \ell\|_0$: số thành phần khác 0 (cardinality).

- $d = 47,236$: tổng số features.
- Sparsity: $\frac{\|\nabla \ell\|_0}{d} \approx 0.14\% \rightarrow 99.86\%$ gradient components = 0.

Ảnh hưởng đến SGD update:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla \ell(\mathbf{w}^{(t)})$$

Khi gradient cực kỳ sparse:

- Chỉ ~ 67 tham số (0.14%) được cập nhật mỗi iteration.
- 47,169 tham số còn lại (99.86%) “đứng yên”.
- Tần suất cập nhật không đồng đều: features phổ biến được cập nhật hàng trăm lần, features hiếm có thể không được cập nhật lần nào.
- Magnitude cập nhật cũng không đồng đều: common features có tổng gradient tích lũy lớn hơn hàng nghìn lần so với rare features.

So sánh với Dense Data:

Đặc điểm	Dense Data	Sparse Data (RCV1)
Features active/sample	$\sim 100\%$	$\sim 0.14\%$
Gradient sparsity	Low ($\sim 0\%$)	Extreme (99.86%)
Update frequency	Đồng đều	Rất không đồng đều
LR cố định	Hoạt động tốt	Thất bại

Kết luận:

- Gradient thưa (sparse gradient) dẫn đến tần suất cập nhật không đồng đều giữa các tham số.
 - SGD với learning rate cố định không thể cân bằng giữa:
 - Features phổ biến: cần LR nhỏ để tránh oscillation.
 - Features hiếm: cần LR lớn để học được trong thời gian hợp lý.
- \Rightarrow Cần cơ chế adaptive learning rate tự động điều chỉnh theo từng tham số.

3 Adagrad (Adaptive Gradient Algorithm) - Giải Pháp

3.1 Động lực: Tốc độ học Thích ứng (Motivation)

- **Ý tưởng:** Cung cấp một learning rate *riêng biệt cho từng tham số* và *tự động điều chỉnh* nó dựa trên "lịch sử" gradient.

3.2 Tiền điều kiện của Adagrad

3.2.1 Giới thiệu

Trong tối ưu hóa, hình dạng của hàm mất mát quyết định tất cả. Với các bài toán Deep Learning, ta thường mô hình hóa hành vi cục bộ bằng **Hàm bậc hai lồi Convex Quadratic Function** :

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x} + b \quad (2)$$

Trên thực tế của Deep Learning , các bài toán là "Non-convex" (Không lồi), có vô số đỉnh, đáy, đèo, vực thẳm. Việc chứng minh toán học chặt chẽ trên hàm này gần như là bất khả thi.

Nên nếu một thuật toán (như Adagrad) hoạt động tốt trên hàm lồi bậc hai (giải quyết được vấn đề méo mó, thung lũng hẹp), thì trực giác cho thấy nó cũng sẽ hoạt động tốt trên các bề mặt phức tạp của Deep Learning.

Lý do Toán học: Mọi thứ đều là bậc hai khi "Zoom" đủ gần (Khai triển Taylor). Đây là lý do quan trọng nhất về mặt kỹ thuật. Theo định lý Taylor, bất kỳ hàm số trơn nào (dù phức tạp đến đâu), nếu bạn phóng to (zoom) vào cực gần một điểm cực tiểu cục bộ, nó đều trông giống hệt một Hàm bậc hai. Công thức khai triển Taylor tại điểm \mathbf{x} gần điểm tối ưu \mathbf{a} :

$$f(\mathbf{x}) \approx f(\mathbf{a}) + \nabla f(\mathbf{a})^\top (\mathbf{x} - \mathbf{a}) + \frac{1}{2}(\mathbf{x} - \mathbf{a})^\top \mathbf{H}(\mathbf{x} - \mathbf{a})$$

Thành phần thứ 3 ($\frac{1}{2}\dots\mathbf{H}\dots$) chính là dạng Quadratic Function đã đề cập (2) ($\frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x}$). Ma trận \mathbf{Q} trong bài học chính là ma trận Hessian (\mathbf{H}) tại điểm đó.

Kết luận: Để hiểu rõ việc tối ưu tại lân cận cho điểm cực tiểu cho bất kỳ bài toán nào, thì cách tốt nhất chính là thông qua học cách tối ưu cho bài toán hàm bậc hai lồi. Nó là hình ảnh đại diện chính xác cho bài toán khi ta gần chạm tới đích (nhờ Taylor).

3.2.2 Sự biến đổi của hàm lỗi bậc hai

Sau khi đã hiểu được lý do chọn hàm lỗi bậc hai để biểu diễn ý nghĩa của việc tối ưu hoá, ta tiếp tục với việc biến đổi hàm ban đầu (2) thành hàm:

$$f(\mathbf{x}) = \bar{f}(\bar{\mathbf{x}}) = \frac{1}{2} \bar{\mathbf{x}}^\top \mathbf{\Lambda} \bar{\mathbf{x}} + \bar{\mathbf{c}}^\top \bar{\mathbf{x}} + b \quad (3)$$

Để chứng minh, ta cần 2 giả định toán học quan trọng :

- Phân rã trị riêng (Eigendecomposition) : $\mathbf{Q} = \mathbf{U}^\top \mathbf{\Lambda} \mathbf{U}$. Với \mathbf{U} là ma trận trực giao (Orthogonal Matrix), có tính chất đặc biệt: $\mathbf{U}^\top = \mathbf{U}^{-1}$ (nghịch đảo bằng chuyển vị) hay $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$. Và $\mathbf{\Lambda}$ là ma trận đường chéo (Diagonal Matrix) chứa các trị riêng.

Giải thích chi tiết: Tại sao $\mathbf{Q} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$?

1. Lý thuyết đằng sau: ĐỊNH LÝ PHỔ (THE SPECTRAL THEOREM)

Trong toán học, định lý này phát biểu rằng: "Mọi ma trận thực đối xứng (Symmetric Real Matrix) đều có thể được chéo hóa bởi một ma trận trực giao."

2. Tại sao áp dụng ở đây?

Ma trận Hessian \mathbf{Q} chứa các đạo hàm riêng bậc hai $\left(\frac{\partial^2 f}{\partial x_i \partial x_j}\right)$. Theo định lý Schwarz trong giải tích, nếu hàm số liên tục thì:

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$$

Điều này nghĩa là \mathbf{Q} luôn đối xứng. Do đó, Định lý Phổ luôn đúng trong trường hợp này.

3. Quy trình tính toán từng bước (Để tìm \mathbf{U} và $\mathbf{\Lambda}$):

1. **Bước 1:** Giải phương trình đặc trưng $\det(\mathbf{Q} - \lambda \mathbf{I}) = 0$.

Kết quả: Tìm được các nghiệm $\lambda_1, \lambda_2, \dots$ (Đây chính là các phần tử nằm trên đường chéo của ma trận $\mathbf{\Lambda}$).

2. **Bước 2:** Với mỗi λ_i , giải hệ phương trình $(\mathbf{Q} - \lambda_i \mathbf{I})\mathbf{x} = 0$.

Kết quả: Tìm được các vector riêng \mathbf{v}_i .

3. **Bước 3 (Quan trọng):** Chuẩn hóa vector riêng (Normalization).

Tính $\mathbf{u}_i = \frac{\mathbf{v}_i}{\|\mathbf{v}_i\|}$ sao cho độ dài vector bằng 1.

Lý thuyết: Bước này đảm bảo ma trận \mathbf{U} là trực giao ($\mathbf{U}^\top = \mathbf{U}^{-1}$), giúp việc xoay trục không làm méo kích thước không gian (bảo toàn độ dài Euclid).

- **Đổi biến số (Change of variables)**: Ta đặt biến mới: $\bar{x} = Ux$. Từ đó suy ra x theo \bar{x} : Nhân hai vế với U^T (nghịch đảo của U), ta được:

$$x = U^T \bar{x}$$

Chúng ta sẽ biến đổi từng thành phần của hàm số ban đầu (2).

Bước A: Biến đổi Thành phần Bậc hai ($\frac{1}{2}x^T Qx$)

Thay thế Q bằng dạng phân rã $U^T \Lambda U$:

$$\text{Term}_1 = \frac{1}{2}x^T (U^T \Lambda U)x$$

Sử dụng tính chất kết hợp của phép nhân ma trận để nhóm lại:

$$\text{Term}_1 = \frac{1}{2}(x^T U^T) \Lambda (Ux)$$

Sử dụng tính chất chuyển vị $(AB)^T = B^T A^T$. Ta thấy cụm $(x^T U^T)$ chính là $(Ux)^T$. Vậy phương trình trở thành:

$$\text{Term}_1 = \frac{1}{2}(Ux)^T \Lambda (Ux)$$

Thay biến mới $\bar{x} = Ux$:

$$\text{Term}_1 = \frac{1}{2}\bar{x}^T \Lambda \bar{x}$$

⇒ Xong phần bậc hai.

Bước B: Biến đổi Thành phần Tuyến tính ($c^T x$)

Thay thế x bằng $U^T \bar{x}$ (như đã rút ra ở mục 2):

$$\text{Term}_2 = c^T (U^T \bar{x})$$

Nhóm các vector hệ số lại:

$$\text{Term}_2 = (c^T U^T) \bar{x}$$

Sử dụng tính chất chuyển vị ngược lại $B^T A^T = (AB)^T$. Cụm $(c^T U^T)$ chính là $(Uc)^T$:

$$\text{Term}_2 = (Uc)^T \bar{x}$$

Tài liệu định nghĩa vector hệ số mới $\bar{c} = U c$. Vậy ta có:

$$\text{Term}_2 = \bar{c}^T \bar{x}$$

⇒ Xong phần tuyến tính.

Bước C: Tổng hợp kết quả

Ghép Bước A và Bước B lại với hằng số b (hằng số b giữ nguyên vì không chứa x), ta có công thức cuối cùng:

$$f(x) \rightarrow \bar{f}(\bar{x}) = \underbrace{\frac{1}{2} \bar{x}^T \Lambda \bar{x}}_{\text{Từ bước A}} + \underbrace{\bar{c}^T \bar{x}}_{\text{Từ bước B}} + b$$

Đây chính xác là công thức 12.7.1 của D2L hay là hàm số (3)

Vấn đề: Nếu ma trận Hessian Q có Số điều kiện (Condition Number) lớn ($\lambda_{max} \gg \lambda_{min}$), hàm số sẽ có dạng thung lũng hẹp

Hậu quả: Gradient Descent thông thường sẽ bị dao động mạnh (zigzag) tại vách núi dựng đứng và di chuyển rùa bò tại lòng thung lũng phẳng.

Vì việc tính toàn bộ trị riêng quá tốn kém ($O(d^3)$), ta dùng phương pháp xấp xỉ: Chuẩn hóa dựa trên đường chéo. Công thức tổng quát:

$$\tilde{Q} = \text{diag}^{-\frac{1}{2}}(Q) Q \text{diag}^{-\frac{1}{2}}(Q) \quad (4)$$

Việc tính các giá trị của ma trận sẽ được tính theo công thức

$$\tilde{Q}_{ij} = \frac{Q_{ij}}{\sqrt{Q_{ii} Q_{jj}}}$$

Quy trình tính toán từng bước (Step-by-step Derivation): Giả sử ta có ma trận Hessian rất "méo":

$$Q = \begin{bmatrix} 100 & 5 \\ 5 & 1 \end{bmatrix}.$$

Bước 1: Tạo ma trận tỉ lệ S

$$\text{Lấy nghịch đảo căn bậc hai đường chéo: } S = \text{diag}\left(\frac{1}{\sqrt{100}}, \frac{1}{\sqrt{1}}\right) = \begin{bmatrix} 0.1 & 0 \\ 0 & 1 \end{bmatrix}$$

Bước 2: Biến đổi ma trận $\tilde{Q} = S Q S$

$$\tilde{Q}_{11} = 100 \times 0.1 \times 0.1 = 1, \tilde{Q}_{22} = 1 \times 1 \times 1 = 1, \tilde{Q}_{12} = 5 \times 0.1 \times 1 = 0.5 \Rightarrow$$

$$\text{Kết quả: } \tilde{Q} = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}.$$

Ý nghĩa: Đường chéo đã được chuẩn hóa về 1. Không gian đã được cân bằng lại, giúp thuật toán di chuyển dễ dàng hơn.

Trong thực tế, ta thậm chí không biết ma trận \mathbf{Q} (Hessian). Làm sao ta áp dụng phép tiền điều kiện trên?

Chúng ta dựa vào Khai triển Taylor bậc 2 quanh điểm cực tiểu \mathbf{x}^* :

$$f(\mathbf{x}) \approx f(\mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^\top \mathbf{H}(\mathbf{x} - \mathbf{x}^*)$$

Lấy đạo hàm (Gradient) hai vế:

$$\mathbf{g} = \nabla f(\mathbf{x}) \approx \mathbf{H}(\mathbf{x} - \mathbf{x}^*)$$

Suy luận Heuristic (Quy tắc ngón tay cái) Từ phương trình trên, nếu ta coi khoảng cách tới đích $(\mathbf{x} - \mathbf{x}^*)$ là đại lượng biến thiên chậm, ta có mối tương quan tỉ lệ thuận:

$$\|\mathbf{g}\| \propto \|\mathbf{H}\|$$

Gradient lớn \Rightarrow Hessian lớn (Độ cong lớn/Đốc đứng). Gradient nhỏ \Rightarrow Hessian nhỏ (Độ cong nhỏ/Phẳng). \Rightarrow Chiến lược: Sử dụng Tổng bình phương Gradient ($\sum \mathbf{g}^2$) làm biến đại diện (Proxy) cho Hessian (\mathbf{Q}) để thực hiện phép chia tỉ lệ.

3.3 Thuật toán Adagrad

- **Biến trạng thái (State Variable):** s_t , một véc-tơ tích lũy *tổng bình phương* của các gradient trong quá khứ.
- **Công thức tích lũy:** (với g_t là gradient tại bước t)

$$s_t = s_{t-1} + g_t^2$$

- **Công thức cập nhật (Update Rule):**

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{s_t + \epsilon}} \odot g_t$$

- η là learning rate **toàn cục**.
- ϵ là hằng số nhỏ (ví dụ $1e-7$) để tránh chia cho 0.

- \odot là phép nhân element-wise.
- Phân số $\frac{\eta}{\sqrt{s_t+\epsilon}}$ là **learning rate thích ứng** cho từng tham số.

Giải thích thêm cho việc tại sao learning rate được chia cho căn bậc 2 tổng tích lũy bình phương

1. Lý do Toán học: Mô phỏng công thức (4)

- Lý thuyết bảo rằng: Muốn chuẩn hóa hình dạng hàm số (biến elip thành hình tròn), ta phải chia tỉ lệ trục tọa độ cho căn bậc hai của độ cong Hessian ($\sqrt{Q_{ii}}$).
- Adagrad : Vì không có Hessian, ta dùng tổng bình phương Gradient (s_t) để thay thế.
- Kết luận: Dấu căn bậc hai $\sqrt{s_t}$ xuất hiện chính là để thực hiện đúng phép toán "mũ $-1/2$ " của lý thuyết

2. Lý do Thống kê: Chuẩn hóa về độ lệch chuẩn (RMS)

Trong thống kê, $s_t = \sum g^2$ đại diện cho tổng phương sai (variance) chưa chuẩn hóa. Nếu Gradient biến động mạnh (lớn), phương sai (s_t) sẽ rất lớn. Để đo "quy mô" (scale) của sự biến động, ta không dùng phương sai, mà dùng Độ lệch chuẩn (Standard Deviation). Công thức độ lệch chuẩn là: $\sigma = \sqrt{\text{Variance}}$. Bằng cách chia cho $\sqrt{s_t}$, Adagrad thực chất đang làm một phép tính gọi là RMS Normalization (Root Mean Square). Nó đưa các gradient có độ lớn khác nhau về cùng một quy mô đơn vị. Gradient lớn \rightarrow Chia cho số lớn. Gradient nhỏ \rightarrow Chia cho số nhỏ. Kết quả là các bước cập nhật trở nên đồng đều hơn.

3. Cơ chế "Phanh"

Hãy tưởng tượng $\sqrt{s_t}$ như một bộ phận cảm biến tốc độ. Tại sao cần Căn bậc 2 mà không để nguyên s_t ?

s_t là tổng của các bình phương (g^2). Nó tăng lên rất nhanh và có giá trị cực lớn. Nếu chia cho s_t (không có căn), mẫu số sẽ quá lớn, làm cho Learning Rate tụt về 0 quá nhanh \rightarrow Mô hình ngừng học ngay lập tức (Early Stopping). Dùng Căn bậc 2 giúp kìm hãm sự tăng trưởng của mẫu số, giữ cho nó tuyến tính với độ lớn của gradient, giúp Learning Rate giảm từ từ một cách hợp lý ($O(t^{-1/2})$) thay vì tụt dốc không phanh ($O(t^{-1})$).

3.4 Ví dụ Tính toán Cụ thể cho Adagrad

Tổng quan: Cơ chế Adagrad và Lý thuyết Khởi tạo Trọng số

Trong phần này, chúng ta sẽ đi sâu vào cách thuật toán Adagrad hoạt động thông qua một ví dụ tính toán thủ công, sau đó giải thích tại sao các con số khởi tạo trong ví dụ lại khác với thực tế triển khai trong các thư viện Deep Learning.

3.4.1 Phần 1: Ví dụ Tính toán "Coordinate-wise" (Từng tọa độ)

Adagrad nổi bật nhờ khả năng điều chỉnh tốc độ học (learning rate) riêng biệt cho từng tham số dựa trên tần suất xuất hiện của chúng. Để minh họa, ta xét bài toán tối ưu hóa với 2 tham số đại diện cho 2 trường hợp đối lập.

1. Thiết lập Bài toán (Setup) Mô hình: 2 trọng số $\mathbf{w} = [w_1, w_2]$.

- w_1 : Đặc trưng phổ biến, xuất hiện nhiều (VD: từ "learning").
- w_2 : Đặc trưng hiếm, ít xuất hiện (VD: từ "preconditioning").

Khởi tạo:

- Trọng số: $\mathbf{w}_0 = [1.0, 1.0]$.
- Biên tích lũy gradient: $\mathbf{s}_0 = [0, 0]$.
- Siêu tham số: Learning Rate $\eta = 0.1$, $\epsilon \approx 0$.

2. Bước 1: Sự cân bằng tự động ($t = 1$) Giả sử tại bước đầu tiên, Gradient trả về cho thấy w_1 có độ dốc lớn gấp 40 lần w_2 .

Input Gradient: $\mathbf{g}_1 = [4.0, 0.1]$.

Quá trình tính toán:

1. **Bình phương Gradient (g^2):** Tính riêng cho từng tọa độ.

$$[4.0^2, 0.1^2] = [16.0, 0.01]$$

2. **Tích lũy vào \mathbf{s} ($\mathbf{s}_t = \mathbf{s}_{t-1} + \mathbf{g}^2$):**

$$\mathbf{s}_1 = [0, 0] + [16.0, 0.01] = [16.0, 0.01]$$

3. **Hệ số điều chỉnh ($\frac{1}{\sqrt{s}}$):** Đây là bước "phanh" hoặc "tăng tốc".

- Với w_1 : $1/\sqrt{16.0} = 0.25$ (Hệ số nhỏ \rightarrow Phanh gấp).
- Với w_2 : $1/\sqrt{0.01} = 10.0$ (Hệ số lớn \rightarrow Tăng tốc).

4. **Cập nhật trọng số:**

$$\Delta w_1 = 0.1 \times 0.25 \times 4.0 = \mathbf{0.1}$$

$$\Delta w_2 = 0.1 \times 10.0 \times 0.1 = \mathbf{0.1}$$

$$\mathbf{w}_1 = [1.0 - 0.1, \quad 1.0 - 0.1] = [0.9, \quad 0.9]$$

Nhận xét: Dù gradient chênh lệch 40 lần, Adagrad đã tự động cân bằng để cả hai tham số di chuyển cùng một quãng đường.

3. Bước 2: Hiệu ứng tích lũy ($t = 2$) Giả sử Gradient tiếp tục duy trì xu hướng cũ: $\mathbf{g}_2 = [2.0, \quad 0.1]$.

Quá trình tính toán:

1. **Tích lũy tiếp vào \mathbf{s} :**

$$\mathbf{s}_2 = \mathbf{s}_1 + \mathbf{g}_2^2 = [16.0, 0.01] + [4.0, 0.01] = [20.0, \quad 0.02]$$

2. **Learning Rate hiệu dụng (Individual Learning Rate):**

- Với w_1 (Dốc/Nhiều): $\eta_{eff} \approx 0.1/\sqrt{20} \approx \mathbf{0.022}$ (Giảm dần).
- Với w_2 (Thoải/Hiếm): $\eta_{eff} \approx 0.1/\sqrt{0.02} \approx \mathbf{0.707}$ (Vẫn rất cao).

3. **Bước nhảy thực tế:**

- w_1 di chuyển một đoạn rất nhỏ: 0.044.
- w_2 di chuyển một đoạn lớn hơn gradient của nó: 0.07.

Kết luận: Mỗi tham số sở hữu một tốc độ học riêng biệt. Tham số "learning" (w_1) học chậm lại để hội tụ, tham số "preconditioning" (w_2) giữ tốc độ cao để bắt kịp.

3.4.2 Phần 2: Lý thuyết Khởi tạo Trọng số (Weight Initialization)

Trong ví dụ trên, ta đã chọn $w_0 = [1.0, 1.0]$. Tuy nhiên, trong thực tế, việc chọn con số này tuân theo các nguyên tắc nghiêm ngặt hơn nhiều.

1. Tại sao trong ví dụ lại chọn $w_0 = [1.0, 1.0]$? Việc chọn giá trị này hoàn toàn phục vụ mục đích Sư phạm (Pedagogical):

- **Tối giản tính toán:** Giúp ta tập trung vào cơ chế của thuật toán thay vì bị rối bởi các con số thập phân phức tạp (như 0.14159...).
- **Cô lập biến số:** Bằng cách để xuất phát điểm ngang bằng nhau, ta chứng minh được rằng sự thay đổi tốc độ học ở Bước 1 và 2 hoàn toàn do Gradient và Adagrad tạo ra, chứ không phải do trọng số ban đầu khác nhau.

2. Tại sao KHÔNG làm thế trong thực tế? (Cạm bẫy Đối xứng) Nếu áp dụng $w_0 = [1.0, 1.0]$ cho code thực tế, bạn sẽ gặp Vấn đề Đối xứng (Symmetry Problem):

- Tất cả nơ-ron thực hiện phép tính y hệt nhau.
- Tất cả Gradient trả về giống hệt nhau.
- Cả mạng nơ-ron khổng lồ sẽ hoạt động như một nơ-ron duy nhất, không thể học được các đặc trưng phức tạp.

3. Giải pháp Thực tế: Bảo toàn Phương sai (Variance Preservation) Trong Deep Learning thực tế (PyTorch/TensorFlow), trọng số được khởi tạo ngẫu nhiên dựa trên nguyên lý: "Phương sai đầu ra phải bằng phương sai đầu vào" để tránh Vanishing/Exploding Gradient.

- **He Initialization (cho ReLU):** Khởi tạo ngẫu nhiên từ phân phối chuẩn với phương sai $\text{Var}(W) = \frac{2}{n_{in}}$. Đây là chuẩn mực hiện đại cho các mạng ConvNet/ResNet.
- **Xavier/Glorot Initialization (cho Sigmoid/Tanh):** Khởi tạo với phương sai $\text{Var}(W) = \frac{1}{n_{in}}$.

Tóm lại

Ví dụ tính toán $w_0 = [1.0, 1.0]$ là một mô hình "trong phòng thí nghiệm" để hiểu cơ chế cân bằng bước nhảy của Adagrad. Còn trong "thế giới thực", ta phải dùng He/Xavier Initialization để phá vỡ tính đối xứng và đảm bảo tín hiệu truyền đi ổn định.

Kết luận ví dụ: Learning rate của w_1 (phổ biến) liên tục giảm, trong khi learning rate của w_2 (hiếm) chỉ giảm khi nó thực sự xuất hiện.

3.5 Phân tích Ưu điểm và Nhược điểm

- **Ưu điểm:**

- Giải quyết hoàn hảo bài toán **gradient thưa** (sparse gradient).
- Tự động điều chỉnh learning rate.

- **Nhược điểm (Vấn đề chí mạng):**

- Biến tích lũy s_t luôn *tăng* và không bao giờ giảm (vì $g_t^2 \geq 0$).
- Dẫn đến learning rate **giảm dần về 0** \implies Thuật toán **dừng học (stop learning)** quá sớm.

3.6 Adagrad trên Python

3.6.1 Minh hoạ tối ưu hoá Adagrad trên hàm số

<https://colab.research.google.com/drive/1InbIefqC95uq7SNLxdy86s0FvWBUjNpm?usp=sharing>

Hàm số dùng để tối ưu hoá là :

$$f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2 \quad (5)$$

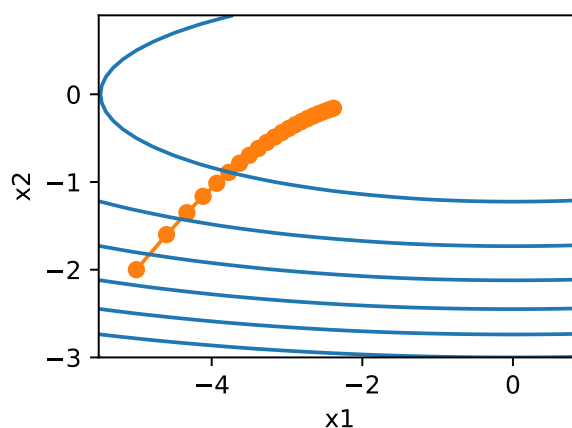
Với learning rate $\eta = 0.4$

```
1 import math
2 import torch
3 from d2l import torch as d2l
4
5 def adagrad_2d(x1, x2, s1, s2):
6     eps = 1e-6
7     g1, g2 = 0.2 * x1, 4 * x2
8     s1 += g1 ** 2
9     s2 += g2 ** 2
10    x1 -= eta / math.sqrt(s1 + eps) * g1
11    x2 -= eta / math.sqrt(s2 + eps) * g2
12    return x1, x2, s1, s2
```

```

13
14 def f_2d(x1, x2):
15     return 0.1 * x1 ** 2 + 2 * x2 ** 2
16
17 eta = 0.4
18 d2l.show_trace_2d(f_2d, d2l.train_2d(adagrad_2d))

```



Hình 7: Minh hoạ đường đi của adagrad

Để hiểu hình này, chúng ta cần phân tích 3 thành phần chính: Địa hình (các đường xanh), Hành trình (đường cam) và Tốc độ (khoảng cách giữa các điểm).

1. Địa hình: Các đường đồng mức màu xanh

- Hình dáng: Đây là một cái bát hình elip bị kéo dãn.
- Trục dọc (x_2): Có hệ số là 2 (lớn). Điều này nghĩa là hướng này rất dốc. Chỉ cần nhích nhẹ x_2 , giá trị hàm số thay đổi rất mạnh. Các đường xanh nằm sát nhau biểu thị độ dốc lớn này.
- Trục ngang (x_1): Có hệ số là 0.1 (nhỏ). Điều này nghĩa là hướng này rất thoải (phẳng). Các đường xanh nằm cách xa nhau.

2. Hành trình: Đường màu cam (Quỹ đạo Adagrad) Các chấm màu cam là vị trí của (x_1, x_2) sau mỗi bước cập nhật (epoch). Đường nối chúng là đường đi của thuật toán.

- Nếu dùng thuật toán thường (như SGD), vì trục dọc (x_2) quá dốc, thuật toán sẽ bị văng qua lại mạnh (zigzag) theo chiều dọc và đi rất chậm theo chiều ngang.

- Adagrad thấy Gradient ở trục dọc (x_2) rất lớn \rightarrow Nó tích lũy vào $s_2 \rightarrow$ và chia learning rate cho số lớn này. Sau đó bước đi theo chiều dọc tự động nhỏ lại (Hãm phanh). Adagrad thấy ở trục ngang (x_1) rất nhỏ \rightarrow nên tích lũy vào s_1 (số nhỏ) \rightarrow và chia learning rate cho số nhỏ \rightarrow . Bước đi theo chiều ngang được giữ nguyên hoặc lớn hơn (Tăng tốc).
- Kết quả: Ta thấy đường màu cam không bị zigzag hỗn loạn mà uốn cong một cách mượt mà về phía tâm. Nó đi cẩn thận ở chỗ nguy hiểm (x_2) và đi dứt khoát ở chỗ an toàn (x_1).

3. Tốc độ: Khoảng cách giữa các chấm cam

Hãy nhìn kỹ các chấm cam:

- Lúc đầu: Các chấm cách xa nhau (bước nhảy lớn).
- Càng về sau: Các chấm xít lại gần nhau dày đặc (bước nhảy li ti).

Lý do: Đây là minh chứng cho việc Learning Rate giảm dần (Decay).

- Công thức cập nhật có mẫu số là $\sqrt{s_t}$.
- Vì s_t cứ cộng dồn mãi ($s_t = s_{t-1} + g^2$), nên mẫu số càng ngày càng to.
- Kết quả là càng về cuối, Adagrad càng "đuối sức", bước đi trở nên rất nhỏ, khiến nó mất nhiều thời gian để nhích từng chút một vào tâm điểm $(0, 0)$.

3.6.2 Xây dựng thuật toán Adagrad từ công thức

Chúng ta sẽ viết code Python từ đầu cho thuật toán Adagrad dựa trên những gì đã nhắc tới và huấn luyện mô hình

```

1
2 d2l.DATA_HUB['airfoil'] = (d2l.DATA_URL + 'airfoil_self_noise.dat',
3                             '76e5be1548fd8222e5074cf0faae75edff8cf93f')
4
5 def get_data_ch11(batch_size=10, n=1500):
6     """Defined in :numref:`sec_minibatches`"""
7     data = np.genfromtxt(d2l.download('airfoil'),
8                           dtype=np.float32, delimiter='\t')
9     data = torch.from_numpy((data - data.mean(axis=0)) /
10                             data.std(axis=0))

```

```
10     data_iter = d2l.load_array((data[:n, :-1], data[:n, -1]),
11                               batch_size, is_train=True)
12     return data_iter, data.shape[1]-1
13
14 def train_ch11(trainer_fn, states, hyperparams,
15               data_iter, feature_dim, num_epochs=2):
16     """Defined in :numref:`sec_minibatches`"""
17     # Initialization
18     w = torch.normal(mean=0.0, std=0.01, size=(feature_dim, 1),
19                   requires_grad=True)
20     b = torch.zeros((1), requires_grad=True)
21     net, loss = lambda X: d2l.linreg(X, w, b), d2l.squared_loss
22     # Train
23     animator = d2l.Animator(xlabel='epoch', ylabel='loss', xlim=[0,
24                               num_epochs], ylim=[0.22, 0.35])
25     n, timer = 0, d2l.Timer()
26     for _ in range(num_epochs):
27         for X, y in data_iter:
28             l = loss(net(X), y).mean()
29             l.backward()
30             trainer_fn([w, b], states, hyperparams)
31             n += X.shape[0]
32             if n % 200 == 0:
33                 timer.stop()
34                 animator.add(n/X.shape[0]/len(data_iter),
35                               (d2l.evaluate_loss(net, data_iter, loss),))
36                 timer.start()
37
38     print(f'loss: {animator.Y[0][-1]:.3f},
39           {timer.sum()/num_epochs:.3f} sec/epoch')
40     return timer.cumsum(), animator.Y[0]
41
42 def init_adagrad_states(feature_dim):
43     s_w = torch.zeros((feature_dim, 1))
44     s_b = torch.zeros(1)
45     return (s_w, s_b)
```

```

42
43 def adagrad(params, states, hyperparams):
44     eps = 1e-6
45     for p, s in zip(params, states):
46         with torch.no_grad():
47             s[:] += torch.square(p.grad)
48             p[:] -= hyperparams['lr'] * p.grad / torch.sqrt(s + eps)
49     p.grad.data.zero_()

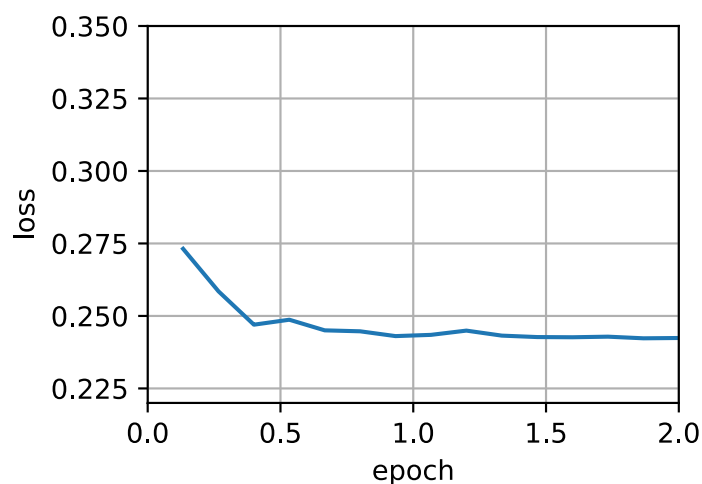
```

Ta thu được kết quả tối ưu loss giống với những gì đã biết

```

1     loss: 0.242, 0.033 sec/epoch

```



Hình 8: Loss của adagrad trên tập airfoil-self-noise sau 20 epochs

3.6.3 Minh họa Adagrad trên tập RCV1

```

1 import torch
2 import torch.nn as nn
3 import time
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from sklearn.datasets import fetch_rcv1
7
8 # 1. Setup & Load Data (Giữ nguyên như code cũ)
9 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

```

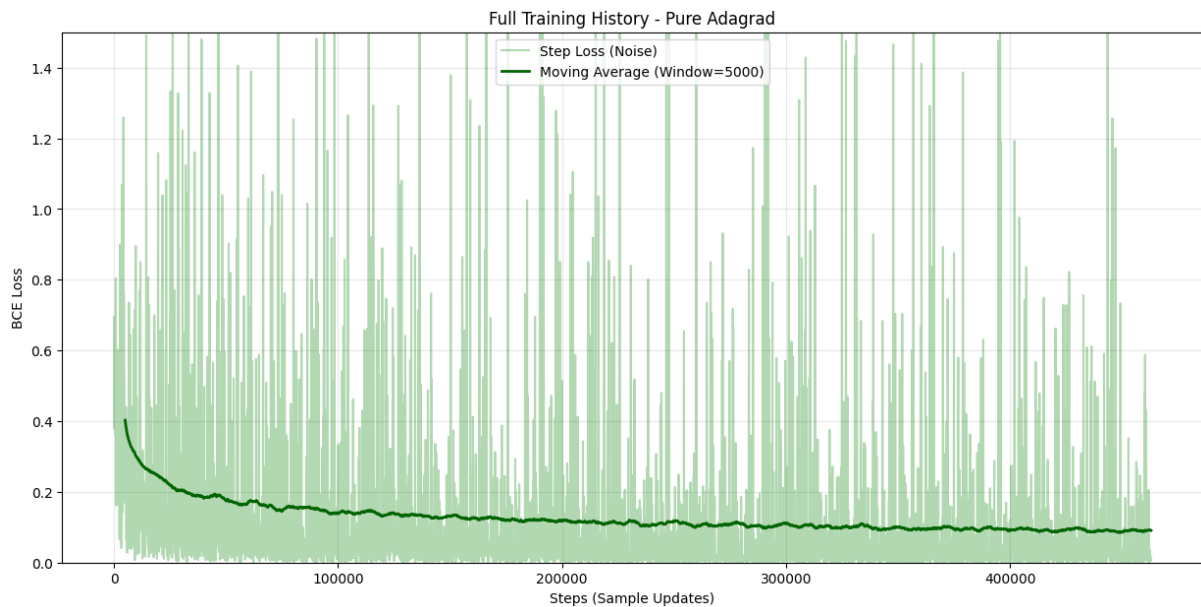


```
10 print(f"Device đang chạy: {device}")
11
12 print("Downloading/Loading RCV1...")
13 rcv1 = fetch_rcv1(subset='train')
14 n_samples, n_features = rcv1.data.shape
15
16 # Xu lý nhãn (Label Processing)
17 topic_counts = rcv1.target.sum(axis=0)
18 most_frequent_topic_idx = np.argmax(topic_counts)
19 y_full = rcv1.target[:, most_frequent_topic_idx].toarray().ravel()
20 y_cpu = torch.tensor(y_full, dtype=torch.float32).view(-1, 1)
21
22 print(f"Dataset Loaded. Samples: {n_samples}, Features:
      {n_features}")
23 print("-" * 50)
24
25 # 2. Hàm Train ADAGRAD (Thay đổi logic update tại đây)
26 def train_adagrad_visualize(X_scipy, y_cpu, lr=0.01, epochs=1,
      epsilon=1e-8):
27     # Khởi tạo tham số W và b
28     W = torch.normal(0, 0.01, size=(n_features, 1),
      requires_grad=True, device=device)
29     b = torch.zeros(1, requires_grad=True, device=device)
30
31     # --- KHOI TẠO BỘ NHỚ ADAGRAD (Accumulated Squared Gradients) ---
32     # G_W và G_b lưu tổng bình phương gradient của W và b
33     G_W = torch.zeros_like(W, device=device)
34     G_b = torch.zeros_like(b, device=device)
35
36     criterion = nn.BCEWithLogitsLoss()
37     step_losses = []
38
39     print(f"Bắt đầu Training Pure ADAGRAD (Total samples:
      {n_samples})...")
40     start_time = time.time()
```

```
41
42     for epoch in range(epochs):
43         indices = np.random.permutation(n_samples)
44
45         for i, idx in enumerate(indices):
46             # Lấy 1 mẫu (Stochastic)
47             X_sample = torch.tensor(X_scipy[idx].toarray(),
48                                     dtype=torch.float32).to(device)
49             y_sample = y_cpu[idx].to(device).view(1, 1)
50
51             # Forward
52             linear_out = X_sample @ W + b
53             loss = criterion(linear_out, y_sample)
54
55             # Backward
56             loss.backward()
57
58             # --- CẬP NHẬT TRỌNG SỐ THEO CÔNG THỨC ADAGRAD ---
59             with torch.no_grad():
60                 # 1. Cộng dồn bình phương gradient vào bộ nhớ ( $G_t = G_{t-1} + g_t^2$ )
61                 G_W += W.grad ** 2
62                 G_b += b.grad ** 2
63
64                 # 2. Tính mẫu số điều chỉnh (căn bậc hai của  $G + \epsilon$ )
65                 std_W = torch.sqrt(G_W) + epsilon
66                 std_b = torch.sqrt(G_b) + epsilon
67
68                 # 3. Update tham số:  $W = W - (lr / std) * grad$ 
69                 # Lưu ý: Mỗi tham số sẽ có learning rate riêng biệt
70                 W -= (lr / std_W) * W.grad
71                 b -= (lr / std_b) * b.grad
72
73                 # 4. Reset gradient
74                 W.grad.zero_()
```

```
74         b.grad.zero_()
75
76         # Lưu lại loss
77         current_loss = loss.item()
78         step_losses.append(current_loss)
79
80         # In tiến độ
81         if i % 5000 == 0:
82             print(f"Epoch {epoch+1} | Sample {i}/{n_samples} | Loss:
83                   {current_loss:.4f}")
84
85             print(f"--> Kết thúc Epoch {epoch+1}")
86
87             print(f"Tổng thời gian: {time.time() - start_time:.2f}s")
88             return step_losses
89
90 # 3. Hàm Visualize (Giữ nguyên để so sánh công bằng)
91 def plot_full_training_history(losses, alg_name="Adagrad",
92                               window_size=1000, downsample_rate=100):
93
94     plt.figure(figsize=(15, 7))
95     total_steps = len(losses)
96
97     # 1. Vẽ Loss thực tế (Downsampling)
98     steps = range(0, total_steps, downsample_rate)
99     sampled_losses = [losses[i] for i in steps]
100
101     plt.plot(steps, sampled_losses, color='green', alpha=0.3,
102             label='Step Loss (Noise)')
103
104     # 2. Vẽ Moving Average
105     if total_steps > window_size:
106         moving_avg = np.convolve(losses,
107                                   np.ones(window_size)/window_size, mode='valid')
108         ma_x_axis = range(window_size-1, total_steps)
```

```
104     plt.plot(ma_x_axis[::downsample_rate],
              moving_avg[::downsample_rate],
105              color='darkgreen', linewidth=2, label=f'Moving Average
              (Window={window_size})')
106
107     plt.title(f'Full Training History - {alg_name}')
108     plt.xlabel('Steps (Sample Updates)')
109     plt.ylabel('BCE Loss')
110     plt.legend()
111     plt.grid(True, alpha=0.3)
112     plt.ylim(0, 1.5)
113     plt.show()
114
115     # --- SỬ DỤNG ---
116     # Với Adagrad, ta thường có thể đặt Learning Rate khởi điểm cao hơn
        SGD một chút
117     # vì nó sẽ tự giảm dần (decay) theo thời gian.
118     # Tuy nhiên, để so sánh, mình giữ lr=0.01 hoặc bạn có thể thử 0.1 để
        thấy nó hội tụ nhanh hơn hẳn.
119
120     LR = 0.1 # Adagrad thích LR khởi đầu lớn hơn SGD
121     EPOCHS = 20
122
123     print(f"Đang chạy Adagrad với LR={LR}...")
124     loss_history_adagrad = train_adagrad_visualize(rcv1.data, y_cpu,
        lr=LR, epochs=EPOCHS)
125
126     plot_full_training_history(loss_history_adagrad, alg_name="Pure
        Adagrad", window_size=5000, downsample_rate=100)
```



Hình 9: Loss của Adagrad sau 20 epoches

```

1   Device đang chạy: cuda
2   Downloading/Loading RCV1...
3   Dataset Loaded. Samples: 23149, Features: 47236
4   -----
5   Đang chạy Adagrad với LR=0.1...
6   Bắt đầu Training Pure ADAGRAD (Total samples: 23149)...
7   Epoch 1 | Sample 0/23149 | Loss: 0.6934
8   Epoch 1 | Sample 5000/23149 | Loss: 0.2486
9   Epoch 1 | Sample 10000/23149 | Loss: 0.0522
10  Epoch 1 | Sample 15000/23149 | Loss: 0.2880
11  Epoch 1 | Sample 20000/23149 | Loss: 0.4600
12  --> Kết thúc Epoch 1
13  Epoch 2 | Sample 0/23149 | Loss: 0.2018
14  Epoch 2 | Sample 5000/23149 | Loss: 0.1191
15  Epoch 2 | Sample 10000/23149 | Loss: 0.0051
16  Epoch 2 | Sample 15000/23149 | Loss: 0.0430
17  Epoch 2 | Sample 20000/23149 | Loss: 0.2349
18  --> Kết thúc Epoch 2
19  Epoch 3 | Sample 0/23149 | Loss: 0.0456
20  Epoch 3 | Sample 5000/23149 | Loss: 0.3661
21  Epoch 3 | Sample 10000/23149 | Loss: 0.0134

```

```
22 Epoch 3 | Sample 15000/23149 | Loss: 0.0023
23 Epoch 3 | Sample 20000/23149 | Loss: 0.1459
24 --> Kết thúc Epoch 3
25 Epoch 4 | Sample 0/23149 | Loss: 0.0720
26 Epoch 4 | Sample 5000/23149 | Loss: 0.5779
27 Epoch 4 | Sample 10000/23149 | Loss: 0.7443
28 Epoch 4 | Sample 15000/23149 | Loss: 0.0098
29 Epoch 4 | Sample 20000/23149 | Loss: 0.2782
30 --> Kết thúc Epoch 4
31 Epoch 5 | Sample 0/23149 | Loss: 0.0288
32 Epoch 5 | Sample 5000/23149 | Loss: 0.0315
33 Epoch 5 | Sample 10000/23149 | Loss: 0.1064
34 Epoch 5 | Sample 15000/23149 | Loss: 0.2290
35 Epoch 5 | Sample 20000/23149 | Loss: 0.8805
36 --> Kết thúc Epoch 5
37 Epoch 6 | Sample 0/23149 | Loss: 0.0606
38 Epoch 6 | Sample 5000/23149 | Loss: 0.0383
39 Epoch 6 | Sample 10000/23149 | Loss: 0.0014
40 Epoch 6 | Sample 15000/23149 | Loss: 0.0263
41 Epoch 6 | Sample 20000/23149 | Loss: 0.1021
42 --> Kết thúc Epoch 6
43 Epoch 7 | Sample 0/23149 | Loss: 0.0100
44 Epoch 7 | Sample 5000/23149 | Loss: 0.0332
45 Epoch 7 | Sample 10000/23149 | Loss: 0.0586
46 Epoch 7 | Sample 15000/23149 | Loss: 0.1028
47 Epoch 7 | Sample 20000/23149 | Loss: 0.0373
48 --> Kết thúc Epoch 7
49 Epoch 8 | Sample 0/23149 | Loss: 0.0097
50 Epoch 8 | Sample 5000/23149 | Loss: 0.1436
51 Epoch 8 | Sample 10000/23149 | Loss: 0.0225
52 Epoch 8 | Sample 15000/23149 | Loss: 0.0790
53 Epoch 8 | Sample 20000/23149 | Loss: 0.0093
54 --> Kết thúc Epoch 8
55 Epoch 9 | Sample 0/23149 | Loss: 0.0555
56 Epoch 9 | Sample 5000/23149 | Loss: 0.0696
```

```
57 Epoch 9 | Sample 10000/23149 | Loss: 0.0428
58 Epoch 9 | Sample 15000/23149 | Loss: 0.0246
59 Epoch 9 | Sample 20000/23149 | Loss: 0.0964
60 --> Kết thúc Epoch 9
61 Epoch 10 | Sample 0/23149 | Loss: 0.0732
62 Epoch 10 | Sample 5000/23149 | Loss: 0.1775
63 Epoch 10 | Sample 10000/23149 | Loss: 0.0130
64 Epoch 10 | Sample 15000/23149 | Loss: 0.0077
65 Epoch 10 | Sample 20000/23149 | Loss: 0.0888
66 --> Kết thúc Epoch 10
67 Epoch 11 | Sample 0/23149 | Loss: 0.0171
68 Epoch 11 | Sample 5000/23149 | Loss: 0.0948
69 Epoch 11 | Sample 10000/23149 | Loss: 0.0053
70 Epoch 11 | Sample 15000/23149 | Loss: 0.0693
71 Epoch 11 | Sample 20000/23149 | Loss: 0.0666
72 --> Kết thúc Epoch 11
73 Epoch 12 | Sample 0/23149 | Loss: 0.1276
74 Epoch 12 | Sample 5000/23149 | Loss: 0.0049
75 Epoch 12 | Sample 10000/23149 | Loss: 0.0249
76 Epoch 12 | Sample 15000/23149 | Loss: 0.0884
77 Epoch 12 | Sample 20000/23149 | Loss: 0.1339
78 --> Kết thúc Epoch 12
79 Epoch 13 | Sample 0/23149 | Loss: 0.0351
80 Epoch 13 | Sample 5000/23149 | Loss: 0.0562
81 Epoch 13 | Sample 10000/23149 | Loss: 0.0437
82 Epoch 13 | Sample 15000/23149 | Loss: 0.0006
83 Epoch 13 | Sample 20000/23149 | Loss: 0.0069
84 --> Kết thúc Epoch 13
85 Epoch 14 | Sample 0/23149 | Loss: 0.0935
86 Epoch 14 | Sample 5000/23149 | Loss: 0.0377
87 Epoch 14 | Sample 10000/23149 | Loss: 0.0021
88 Epoch 14 | Sample 15000/23149 | Loss: 0.0084
89 Epoch 14 | Sample 20000/23149 | Loss: 0.5105
90 --> Kết thúc Epoch 14
91 Epoch 15 | Sample 0/23149 | Loss: 0.0062
```



```
92 Epoch 15 | Sample 5000/23149 | Loss: 0.0066
93 Epoch 15 | Sample 10000/23149 | Loss: 0.0436
94 Epoch 15 | Sample 15000/23149 | Loss: 0.4036
95 Epoch 15 | Sample 20000/23149 | Loss: 0.0064
96 --> Kết thúc Epoch 15
97 Epoch 16 | Sample 0/23149 | Loss: 0.0184
98 Epoch 16 | Sample 5000/23149 | Loss: 0.0888
99 Epoch 16 | Sample 10000/23149 | Loss: 0.0052
100 Epoch 16 | Sample 15000/23149 | Loss: 0.8732
101 Epoch 16 | Sample 20000/23149 | Loss: 0.2659
102 --> Kết thúc Epoch 16
103 Epoch 17 | Sample 0/23149 | Loss: 0.2450
104 Epoch 17 | Sample 5000/23149 | Loss: 0.0006
105 Epoch 17 | Sample 10000/23149 | Loss: 0.0384
106 Epoch 17 | Sample 15000/23149 | Loss: 0.0063
107 Epoch 17 | Sample 20000/23149 | Loss: 0.3692
108 --> Kết thúc Epoch 17
109 Epoch 18 | Sample 0/23149 | Loss: 0.0658
110 Epoch 18 | Sample 5000/23149 | Loss: 0.0145
111 Epoch 18 | Sample 10000/23149 | Loss: 0.1585
112 Epoch 18 | Sample 15000/23149 | Loss: 0.2380
113 Epoch 18 | Sample 20000/23149 | Loss: 0.0102
114 --> Kết thúc Epoch 18
115 Epoch 19 | Sample 0/23149 | Loss: 0.0852
116 Epoch 19 | Sample 5000/23149 | Loss: 0.0863
117 Epoch 19 | Sample 10000/23149 | Loss: 0.3230
118 Epoch 19 | Sample 15000/23149 | Loss: 0.0170
119 Epoch 19 | Sample 20000/23149 | Loss: 0.0589
120 --> Kết thúc Epoch 19
121 Epoch 20 | Sample 0/23149 | Loss: 0.0017
122 Epoch 20 | Sample 5000/23149 | Loss: 0.1701
123 Epoch 20 | Sample 10000/23149 | Loss: 0.0348
124 Epoch 20 | Sample 15000/23149 | Loss: 0.0417
125 Epoch 20 | Sample 20000/23149 | Loss: 0.2488
126 --> Kết thúc Epoch 20
```


3.7 So sánh giữa Adagrad và SGD trên tập RCV1

3.7.1 Thiết lập Thực nghiệm

- **Tập dữ liệu:** RCV1 (Reuters Corpus Volume I).
- **Số lượng mẫu:** 23,149.
- **Số lượng đặc trưng:** 47,236 (Dữ liệu cao chiều và thưa).
- **Môi trường:** CUDA (GPU RTX3080).
- **Tham số:** Cả hai thuật toán được khởi chạy với cùng thiết lập learning rate ban đầu để đảm bảo tính công bằng trong so sánh.

3.7.2 Kết quả Thực nghiệm

Dựa trên log huấn luyện qua 20 epochs, chúng tôi tổng hợp các chỉ số quan trọng tại các mốc thời gian đại diện trong Bảng 2.

Bảng 2: So sánh Loss và Độ ổn định giữa Adagrad và SGD qua các Epoch

Metric	Epoch	Adagrad Loss	SGD Loss	Nhận xét
Khởi tạo	1	0.6934	0.6814	Tương đương
Hội tụ sớm	3	0.0023	0.5129	Adagrad hội tụ cực nhanh
Giữa kỳ	12	0.1339	1.0979	SGD bị văng (Spike)
Ổn định	16	0.2659	0.1450	Cả hai đều dao động
Kết thúc	20	0.0017	0.2846	Adagrad đạt tối ưu sâu hơn
Thời gian	Total	538.94s	496.02s	SGD nhanh hơn $\approx 43s$

3.7.3 So sánh giữa hai thuật toán tối ưu

Khả năng xử lý dữ liệu thưa (Sparsity Handling) Tập dữ liệu RCV1 có đặc tính là các đặc trưng (từ vựng) xuất hiện với tần suất rất khác nhau.

- **SGD:** Sử dụng learning rate cố định cho toàn bộ 47,236 đặc trưng. Điều này dẫn đến việc các đặc trưng hiếm (rare features) không được cập nhật đủ lớn, trong khi các đặc trưng

phổ biến lại bị cập nhật quá đà. Minh chứng rõ nhất là tại **Epoch 12**, SGD có loss vọt lên 1.0979, cho thấy hiện tượng "overshooting" (văng khỏi đáy).

- **Adagrad:** Cơ chế $\frac{\eta}{\sqrt{s_t}}$ cho phép thuật toán tự động gán learning rate lớn cho các tham số có gradient thưa thớt. Kết quả là ngay tại **Epoch 3**, Loss của Adagrad đã giảm sâu xuống 0.0023, chứng tỏ khả năng "bắt" các đặc trưng quan trọng cực kỳ hiệu quả.

Độ ổn định hội tụ (Convergence Stability) Quan sát hành vi loss ở 5 epoch cuối (16-20):

- **SGD:** Vẫn thể hiện sự bất ổn định cao với các giá trị loss dao động biên độ lớn (từ 0.06 đến 0.93 tại Epoch 18). Điều này cho thấy thuật toán vẫn đang dao động quanh vùng tối ưu mà không thể hội tụ chặt chẽ.
- **Adagrad:** Duy trì mức loss thấp và ổn định hơn. Mặc dù vẫn có dao động (do tính chất Stochastic), nhưng biên độ nhỏ hơn đáng kể. Điều này phù hợp với lý thuyết: learning rate của Adagrad giảm dần theo thời gian ($O(t^{-1/2})$), đóng vai trò như cơ chế "hãm phanh" (annealing) giúp thuật toán tinh chỉnh (fine-tune) chính xác vào điểm cực tiểu.

Chi phí tính toán (Computational Cost) Về mặt thời gian thực thi:

$$T_{Adagrad}(538.94s) > T_{SGD}(496.02s)$$

Sự chênh lệch này là khoảng 8.6%. Nguyên nhân đến từ việc Adagrad phải duy trì và cập nhật thêm ma trận tích lũy gradient s_t và thực hiện phép tính căn bậc hai tại mỗi bước lặp. Tuy nhiên, xét trên hiệu quả hội tụ (đạt loss thấp nhanh gấp 4-5 lần), sự đánh đổi này là hoàn toàn xứng đáng.

3.8 Bài tập – Chương 11.7 Adagrad (D2L)

Ghi chú

Ký hiệu: $\|\cdot\|_2$ là chuẩn Euclid; I là ma trận đơn vị; Q^\top là chuyển vị của Q .

Bài 1

Đề bài. Cho ma trận trực giao \mathbf{U} và vector \mathbf{c}, δ . Chứng minh

$$\|\mathbf{c} - \delta\|_2 = \|\mathbf{U}\mathbf{c} - \mathbf{U}\delta\|_2.$$

Tại sao điều này có nghĩa là độ lớn của nhiễu loạn không thay đổi sau khi thay đổi trực giao của các biến?

Lời giải. Đặt $\mathbf{v} = \mathbf{c} - \delta$. Khi đó cần chứng minh

$$\|\mathbf{v}\|_2 = \|\mathbf{U}\mathbf{v}\|_2.$$

Ta có

$$\|\mathbf{U}\mathbf{v}\|_2^2 = (\mathbf{U}\mathbf{v})^\top (\mathbf{U}\mathbf{v}) = \mathbf{v}^\top \mathbf{U}^\top \mathbf{U} \mathbf{v}.$$

Vì \mathbf{U} là ma trận trực giao nên $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$. Do đó

$$\|\mathbf{U}\mathbf{v}\|_2^2 = \mathbf{v}^\top \mathbf{v} = \|\mathbf{v}\|_2^2.$$

Lấy căn bậc hai hai vế (cả hai không âm) thu được kết quả mong muốn:

$$\|\mathbf{U}\mathbf{v}\|_2 = \|\mathbf{v}\|_2.$$

Giải thích ý nghĩa. Ma trận trực giao là một isometry (phép giữ khoảng cách) trong không gian Euclid: nó biểu diễn các phép quay hoặc phản chiếu (hoặc tổ hợp của chúng) nên không thay đổi độ dài vector và góc giữa các vector. Vì vậy độ lớn của nhiễu (perturbation) $\mathbf{c} - \delta$ không đổi khi ta thay đổi cơ sở bằng một biến đổi trực giao \mathbf{U} .

Bài 2

Đề bài. Thử Adagrad với

$$f_1(\mathbf{x}) = 0.1 x_1^2 + 2 x_2^2$$

và với hàm đã xoay 45 độ

$$f_2(\mathbf{x}) = 0.1 (x_1 + x_2)^2 + 2 (x_1 - x_2)^2.$$

So sánh hành vi của Adagrad trên hai hàm này và giải thích.

1. Phân tích lý thuyết

1. Hessians và condition number.

$$\nabla^2 f_1 = \begin{pmatrix} 0.2 & 0 \\ 0 & 4 \end{pmatrix}, \quad \lambda_{\max}(\nabla^2 f_1) = 4, \quad \lambda_{\min}(\nabla^2 f_1) = 0.2,$$

nên condition number $\kappa = 4/0.2 = 20$.

Vì f_2 là f_1 sau một phép quay 45° (cụ thể tồn tại ma trận trực giao Q sao cho $f_2(\mathbf{x}) = f_1(Q\mathbf{x})$), nên Hessian của f_2 có cùng phổ (cùng các giá trị riêng) như của f_1 . Về mặt phổ, hai hàm có cùng độ “khó” (cùng κ).

2. Tính chất của Adagrad. Adagrad duy trì cho mỗi tọa độ i một accumulator:

$$s_{t,i} = \sum_{\tau=1}^t g_{\tau,i}^2,$$

và cập nhật theo:

$$x_{t+1,i} = x_{t,i} - \eta \frac{g_{t,i}}{\sqrt{s_{t,i}} + \varepsilon}.$$

Đây là phương pháp *coordinate-wise adaptive*: luật điều chỉnh learning-rate phụ thuộc vào lịch sử gradient tại từng tọa độ. Do đó Adagrad **không** bất biến dưới phép quay không gian tham số — nghĩa là nếu ta quay hệ tọa độ, ghi lại gradient theo các tọa độ mới, accumulator per-coordinate sẽ khác và cách cập nhật cũng khác.

3. Hệ quả.

- Với f_1 (axis-aligned): gradient ở mỗi bước thường tập trung vào một tọa độ (do Hessian là đường chéo) — accumulator cho từng tọa độ sẽ phản ánh chính xác mức “stiffness” theo từng tọa độ, nên Adagrad có khả năng điều chỉnh tốt (giảm mạnh lr ở tọa độ có gradient lớn).
- Với f_2 (xoay 45°): gradient của hàm tại mỗi bước thường phân phối qua cả hai tọa độ; do đó cả hai accumulator sẽ tăng theo cách khác, và sự điều chỉnh learning-rate theo tọa độ có thể kém “điểm rơi” hơn so với trường hợp trục-điều hướng, dẫn tới đường đi tối ưu hóa khác (và có thể chậm hoặc dao động hơn).

2. Minh họa bằng số (ý tưởng thí nghiệm)

- **Thiết lập:** chọn cùng điểm khởi tạo $\mathbf{x}^{(0)}$, cùng η (ví dụ 0.1), cùng ε (ví dụ $1e-8$), cùng số bước T .
- **Theo dõi:** giá trị hàm $f(\mathbf{x}_t)$ theo bước, $\|\mathbf{x}_t\|$, và các giá trị $s_{t,i}$ (accumulators).
- **Kỳ vọng quan sát:** trên f_1 Adagrad có thể giảm f nhanh hơn ban đầu theo một số trục; trên f_2 có thể thấy đường đi zig-zag khác và độ suy giảm chậm hơn.

3. Mã Python minh họa (so sánh Adagrad trên f_1 và f_2)

```
import numpy as np
import matplotlib.pyplot as plt

def adagrad_step(x, grad, s, lr=0.1, eps=1e-8):
    s += grad**2
    x = x - lr * grad / (np.sqrt(s) + eps)
    return x, s

def f1(x):
    return 0.1*x[0]**2 + 2*x[1]**2

def grad_f1(x):
    return np.array([0.2*x[0], 4.0*x[1]])

def f2(x):
    # equivalent to f1(Q x) where Q is 45 degree rotation
    # but we can compute gradient by chain rule
    # f2(x) = 0.1*(x1+x2)^2 + 2*(x1-x2)^2
    g1 = 0.2*(x[0]+x[1]) + 4*(x[0]-x[1])
    g2 = 0.2*(x[0]+x[1]) - 4*(x[0]-x[1])
    return np.array([g1, g2])

# init
x0 = np.array([3.0, 1.0])
```

```
lr = 0.2
T = 200
eps = 1e-8

# run on f1
x = x0.copy()
s = np.zeros(2)
hist_f1 = []
for t in range(T):
    g = grad_f1(x)
    x, s = adagrad_step(x, g, s, lr, eps)
    hist_f1.append(0.1*x[0]**2 + 2*x[1]**2)

# run on f2
x = x0.copy()
s = np.zeros(2)
hist_f2 = []
for t in range(T):
    g = f2(x)
    x, s = adagrad_step(x, g, s, lr, eps)
    hist_f2.append(0.1*(x[0]+x[1])**2 + 2*(x[0]-x[1])**2)

plt.semilogy(hist_f1, label='f1 (axis aligned)')
plt.semilogy(hist_f2, label='f2 (rotated)')
plt.xlabel('step')
plt.ylabel('f(x)')
plt.legend()
plt.show()
```

4. Diễn giải kết quả

- Nếu đồ thị cho thấy f_1 giảm nhanh hơn: đó là hệ quả của tính coordinate-adaptive của Adagrad (hàm axis-aligned đồng bộ hơn với update per-coordinate).
- Nếu hai đồ thị tương tự: điều đó có thể xảy ra với các siêu tham số khác nhau hoặc điểm

khởi tạo đặc biệt; tuy nhiên nói chung Adagrad không invariant dưới quay.

Bài 3

Đề bài. Chứng minh định lý vòng tròn Gerschgorin: với ma trận $\mathbf{M} = (m_{ij}) \in \mathbb{C}^{n \times n}$, mọi giá trị riêng λ của \mathbf{M} thỏa mãn

$$|\lambda - m_{jj}| \leq \sum_{k \neq j} |m_{jk}|$$

cho ít nhất một chỉ số j .

Chứng minh. Giả sử λ là eigenvalue của \mathbf{M} với eigenvector tương ứng $\mathbf{x} = (x_1, \dots, x_n)^\top \neq \mathbf{0}$. Chọn j sao cho $|x_j| = \max_{1 \leq i \leq n} |x_i| > 0$. Từ phương trình $\mathbf{M}\mathbf{x} = \lambda\mathbf{x}$, lấy thành phần j ta có

$$(\mathbf{M}\mathbf{x})_j = (\lambda\mathbf{x})_j \implies \sum_{k=1}^n m_{jk}x_k = \lambda x_j.$$

Chuyển số hạng $m_{jj}x_j$ sang về trái:

$$\sum_{k \neq j} m_{jk}x_k = (\lambda - m_{jj})x_j.$$

Lấy trị tuyệt đối cả hai vế và áp dụng bất đẳng thức tam giác:

$$|\lambda - m_{jj}| |x_j| = \left| \sum_{k \neq j} m_{jk}x_k \right| \leq \sum_{k \neq j} |m_{jk}| |x_k|.$$

Vì $|x_k| \leq |x_j|$ theo cách chọn j , suy ra

$$|\lambda - m_{jj}| |x_j| \leq \sum_{k \neq j} |m_{jk}| |x_j| = \left(\sum_{k \neq j} |m_{jk}| \right) |x_j|.$$

Chia cả hai vế cho $|x_j| > 0$ được

$$|\lambda - m_{jj}| \leq \sum_{k \neq j} |m_{jk}|.$$

Đó là điều phải chứng minh. □

Bài 4

Đề bài. Định lý Gerschgorin cho chúng ta biết gì về các trị riêng của ma trận tiền điều hòa theo đường chéo

$$\text{diag}^{-1/2}(\mathbf{M}) \mathbf{M} \text{diag}^{-1/2}(\mathbf{M})?$$

Phân tích. Gọi $\mathbf{D} = \text{diag}(\mathbf{M})$ (ma trận đường chéo chứa các phần tử m_{ii}) và xét ma trận tiền điều hòa

$$\mathbf{P} = \mathbf{D}^{-1/2} \mathbf{M} \mathbf{D}^{-1/2}.$$

Thành phần $P_{ij} = \frac{m_{ij}}{\sqrt{m_{ii}m_{jj}}}$. Do đó:

$$P_{ii} = 1, \quad \text{và} \quad P_{ij} = \frac{m_{ij}}{\sqrt{m_{ii}m_{jj}}} \quad (i \neq j).$$

Áp dụng Gerschgorin cho \mathbf{P} : mọi eigenvalue λ của \mathbf{P} nằm trong ít nhất một đĩa

$$|\lambda - 1| \leq \sum_{j \neq i} \left| \frac{m_{ij}}{\sqrt{m_{ii}m_{jj}}} \right| \quad \text{với một số } i.$$

Nói cách khác, mọi trị riêng của ma trận đã tiền điều hòa tập trung quanh 1, với bán kính bằng tổng các phần tử ngoài đường chéo đã được chuẩn hóa theo $\sqrt{m_{ii}m_{jj}}$. Điều này cho thấy nếu các phần tử ngoài đường chéo m_{ij} nhỏ so với $\sqrt{m_{ii}m_{jj}}$, thì tất cả trị riêng của \mathbf{P} nằm gần 1 (tức ma trận gần \mathbf{I}) — điều này tương ứng với một ma trận được "cân bằng" tốt hơn, giúp cải thiện điều kiện (conditioning) cho các phương pháp lặp và cho tối ưu hóa tuyến tính.

Bài 5

Đề bài. Thử Adagrad cho một mạng sâu (ví dụ Fashion MNIST). Mô tả những gì nên thử và kỳ vọng.

1. Mục tiêu thí nghiệm Đánh giá hiệu năng của Adagrad trên bài toán phân loại hình ảnh (Fashion MNIST) với mạng CNN cơ bản, so sánh với các bộ tối ưu khác (SGD, SGD+momentum, RMSProp, Adam). Quan tâm tới:

- tốc độ hội tụ (train loss giảm theo epoch),
- khả năng tổng quát hoá (validation accuracy),

- hành vi learning-rate hiệu dụng per-parameter (accumulators),
- và tính nhạy với siêu tham số (lr ban đầu, eps).

2. Chi tiết kiến trúc và chuẩn bị

- **Dữ liệu:** Fashion MNIST (60k train, 10k test) với chuẩn hoá đơn giản (div 255, có thể z-score).

- **Kiến trúc ví dụ (gợi ý):**

`Conv(32, 3x3) -> ReLU -> MaxPool(2x2)`

`Conv(64, 3x3) -> ReLU -> MaxPool(2x2)`

`Flatten -> FC(128) -> ReLU -> FC(10) -> Softmax`

- **Siêu tham số đề xuất:**

- batch size: 64,
- epochs: 20–50,
- lr cho Adagrad: 0.01, 0.05, 0.1 (thử nhiều giá trị),
- eps: 1e-8,
- khởi tạo seed cố định để so sánh công bằng.

3. Các thí nghiệm cần làm (chi tiết)

1. **So sánh bộ tối ưu:** chạy với Adagrad, SGD, SGD+momentum (momentum=0.9), RMSProp, Adam. Ghi lại train loss, val loss, val acc theo epoch.
2. **Nhạy cảm với lr:** cho Adagrad thử vài lr khởi điểm để thấy sensitivity; có thể thấy lr lớn dẫn đến quá nhảy, lr nhỏ dẫn đến chậm.
3. **Đo accumulator per-parameter:** (chỉ mẫu nhỏ) in ra/plot distribution của $\sqrt{s_{t,i}}$ cho từng tầng để xem parameter nào bị giảm lr nhiều hơn.
4. **Kiểm tra overfitting/underfitting:** so sánh learning curves.
5. **Thử Adagrad với và không có data augmentation:** xem ảnh hưởng tới generalization.

4. Kỳ vọng (dựa trên hiểu biết về Adagrad)

- **Ưu điểm:** Adagrad tự động giảm lr cho tham số có gradient lớn (thường ổn với dữ liệu sparse hoặc đặc trưng hiếm).
- **Nhược điểm:** accumulator cộng dồn vô hạn có thể làm learning-rate per-parameter giảm quá mức (vanishing step sizes) → training dừng cải thiện sớm; trong deep learning điều này thường làm Adagrad hoạt động *tạm ổn ban đầu* nhưng *thua Adam / RMSProp* về kết quả cuối cùng khi không điều chỉnh kỹ lr.
- **Đo đó:** có khả năng thấy Adagrad hội tụ ban đầu, nhưng bị chững lại sau vài epoch; biến thể như RMSProp/Adam thường cho tốc độ và kết quả tốt hơn.

5. Mã Python (PyTorch) minh họa

Mã mẫu (PyTorch) - minh họa luồng thử nghiệm

```
import torch
```

```
import torch.nn as nn
```

```
import torch.optim as optim
```

```
from torchvision import datasets, transforms
```

```
from torch.utils.data import DataLoader
```

Dữ liệu

```
transform = transforms.Compose([transforms.ToTensor()])
```

```
trainset = datasets.FashionMNIST(root='./data', train=True, download=True,
```

```
testset = datasets.FashionMNIST(root='./data', train=False, download=True)
```

```
train_loader = DataLoader(trainset, batch_size=64, shuffle=True)
```

```
test_loader = DataLoader(testset, batch_size=256, shuffle=False)
```

Mô hình đơn giản

```
class SimpleCNN(nn.Module):
```

```
    def __init__(self):
```

```
        super().__init__()
```

```
        self.conv = nn.Sequential(
```

```
            nn.Conv2d(1, 32, 3, padding=1), nn.ReLU(),
```

```
            nn.MaxPool2d(2),
```

```
        nn.Conv2d(32, 64, 3, padding=1), nn.ReLU(),
        nn.MaxPool2d(2)
    )
    self.fc = nn.Sequential(
        nn.Flatten(),
        nn.Linear(64*7*7, 128), nn.ReLU(),
        nn.Linear(128, 10)
    )
    def forward(self, x): return self.fc(self.conv(x))

device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model = SimpleCNN().to(device)

# Optimizer: Adagrad
optimizer = optim.Adagrad(model.parameters(), lr=0.05, eps=1e-8)
criterion = nn.CrossEntropyLoss()

# Training loop (sơ bộ)
def train_epoch():
    model.train()
    total_loss = 0.0
    for x, y in train_loader:
        x, y = x.to(device), y.to(device)
        optimizer.zero_grad()
        out = model(x)
        loss = criterion(out, y)
        loss.backward()
        optimizer.step()
        total_loss += loss.item() * x.size(0)
    return total_loss / len(train_loader.dataset)

def eval():
    model.eval()
```

```
correct = 0
total = 0

with torch.no_grad():
    for x,y in test_loader:
        x,y = x.to(device), y.to(device)
        out = model(x)
        pred = out.argmax(dim=1)
        correct += (pred==y).sum().item()
        total += y.size(0)

return correct/total

for epoch in range(1,21):
    tr_loss = train_epoch()
    val_acc = eval()
    print(f'Epoch {epoch}: train_loss={tr_loss:.4f}, val_acc={val_acc:.4f}')
```

Bài 6

Đề bài. Cần sửa đổi Adagrad như thế nào để đạt được sự phân rã (decay) ít hung hăng hơn trong learning rate?

1. Vấn đề của Adagrad gốc Adagrad cộng dồn bình phương gradient vô hạn:

$$s_{t,i} = \sum_{\tau=1}^t g_{\tau,i}^2.$$

Do đó $s_{t,i}$ tăng không giảm (non-decreasing) và có thể trở nên rất lớn \rightarrow thành phần $\eta/\sqrt{s_{t,i}}$ giảm rất mạnh theo thời gian, có thể làm bước cập nhật gần bằng 0 trước khi tối ưu tìm tới cực tiểu tốt. Đây là “quá trình phân rã hung hăng” mà đề bài nói.

2. Các sửa đổi phổ biến (chi tiết)

1. RMSProp (EMA của bình phương gradient).

$$s_t = \rho s_{t-1} + (1 - \rho) g_t^2, \quad \rho \in [0.9, 0.99].$$

Thay vì cộng dồn vô hạn, ta sử dụng trung bình mũ (exponential moving average) nên chỉ nhớ lịch sử gần. Điều này ngăn s_t tăng lên vô hạn và giúp learning-rate per-parameter ổn định hơn.

2. **Adadelta**. - Ý tưởng: loại bỏ dependency vào learning-rate toàn cục bằng cách sử dụng tỉ số giữa RMS của update và RMS của gradient:

$$\Delta\theta_t = -\frac{\text{RMS}(\Delta\theta)_{t-1}}{\text{RMS}(g)_t} g_t.$$

- Adadelta dùng EMA cả cho gradient và cho update, tránh giảm lr quá mức.

3. **Adam (với bias-correction)**.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t, \quad v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2,$$

sau đó dùng dạng bias-corrected $m_t/(1 - \beta_1^t)$ và $v_t/(1 - \beta_2^t)$. Adam dùng EMA cho bình phương gradient (như RMSProp) và có moment, dẫn tới ổn định hơn.

4. **AdaBound**. - Giới hạn learning-rate per-parameter trong một khoảng thay đổi theo thời gian để tránh quá lớn hoặc quá nhỏ.
5. **Tăng ε hoặc sử dụng smoothing**. - Dùng ε lớn hơn (ví dụ $1e-6$ thay vì $1e-8$) làm cho denominator không quá lớn ban đầu.
6. **Sử dụng restart hoặc reset accumulator theo epoch**. - thỉnh thoảng reset $s_{t,i}$ về 0 (hoặc nhân với factor nhỏ) để tránh bị đóng băng.
7. **Sử dụng learning-rate schedule toàn cục kết hợp**. - Dùng lr schedule (decay chậm) trên lr ban đầu thay vì phụ thuộc hoàn toàn vào accumulator.

3. Khi nào nên dùng sửa đổi nào?

- Nếu bạn muốn sửa nhanh: chuyển sang RMSProp hoặc Adam — đây là lựa chọn thực tế và hiệu quả.
- Nếu bài toán *sparse* (ví dụ NLP với features hiếm): Adagrad gốc có thể vẫn tốt; nếu bị giảm quá mạnh, thử tăng ε hoặc reset accumulator.
- Nếu cần kiểm soát chặt chẽ: dùng AdaBound hoặc clipping upper/lower bound cho lr per-parameter.

4. Mã ví dụ (RMSProp nhỏ thay cho Adagrad)

```
# RMSProp replacement for Adagrad (numpy style simple)
import numpy as np

x = np.array([3.0, 1.0])
s = np.zeros_like(x)
lr = 0.01
rho = 0.9
eps = 1e-8

for t in range(100):
    g = grad_f1(x)    # giả sử grad_f1 được định nghĩa
    s = rho * s + (1 - rho) * (g**2)
    x = x - lr * g / (np.sqrt(s) + eps)
```

5. Kết luận Để giảm tính “hung hăng” của việc giảm learning-rate trong Adagrad, cách thực tế và phổ biến nhất là dùng **EMA của bình phương gradient** (RMSProp/Adam) hoặc dùng các biến thể như Adadelta/AdaBound; nếu không muốn đổi thuật toán thì có thể chỉnh ε , reset accumulator, hoặc áp một schedule nhân học cho learning-rate toàn cục.

4 Tổng kết và Hướng phát triển

4.1 Tóm tắt

- Tóm tắt lại hành trình:
 1. Bài toán: Tối ưu hóa $L(W)$.
 2. Thách thức: Local Minima, Saddle Points, Vanishing Gradient.
 3. Giải pháp cơ bản: Dùng SGD ∇L .
 4. Hạn chế của SGD: Vấn đề **Gradient Thưa**.
 5. Giải pháp thích ứng: Adagrad.

4.2 Các hướng cải tiến (Future Work)

- Chỉ ra nhược điểm "dừng học sớm" của Adagrad chính là động lực cho các thuật toán tiếp theo.
 - **RMSProp**: Giải quyết vấn đề s_t tăng mãi bằng cách sử dụng *trung bình trượt theo hàm mũ (exponential moving average)*.
 - **Adam**: Kết hợp ý tưởng của RMSProp và Momentum.

Tài liệu

- [1] Duchi, J., Hazan, E., & Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul), 2121-2159.
- [2] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [3] Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2023). *Dive into Deep Learning*. Cambridge University Press.
- [4] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li, “RCV1: A new benchmark collection for text categorization research,” *Journal of Machine Learning Research*, vol. 5, pp. 361–397, 2004.