

# I. Các phương pháp Optimizer:

## 1. Giới thiệu:

Mô hình học máy ngày càng trở nên phổ biến trong nhiều lĩnh vực do khả năng học tập và tự điều chỉnh. Tuy nhiên, quá trình huấn luyện mô hình đòi hỏi sự cân nhắc kỹ lưỡng về các phương pháp tối ưu hóa (Optimizer) để tối ưu hóa hàm mất mát. Trong nghiên cứu này, chúng ta sẽ tìm hiểu và so sánh các phương pháp Optimizer phổ biến để đánh giá ảnh hưởng của chúng đối với quá trình huấn luyện. Dưới đây là các phương pháp Optimizer trong huấn luyện mô hình học máy:

## 2. Stochastic Gradient (SGD):

Stochastic Gradient Descent (SGD) là một trong những phương pháp tối ưu hóa cơ bản và phổ biến nhất được sử dụng trong quá trình huấn luyện mô hình học máy. Dưới đây là giải thích chi tiết về cách SGD hoạt động:

- Tính Gradient: Cho mỗi mini-batch (một lượng nhỏ mẫu dữ liệu được chọn ngẫu nhiên từ tập huấn luyện), tính gradient của hàm mất mát đối với từng trọng số của mô hình.
- Cập Nhật Trọng Số: Cập nhật trọng số của mô hình theo hướng ngược của gradient
- Lặp Lại Quá Trình: Lặp lại quá trình trên với các mini-batch khác nhau cho đến khi đã sử dụng hết toàn bộ dữ liệu huấn luyện (một epoch). Quá trình lặp này có thể được thực hiện qua nhiều epoch để cải thiện hiệu suất của mô hình.
- Lựa Chọn Tỷ Lệ Học ( $\eta$ ): Quá Lớn: Có thể làm cho mô hình không hội tụ hoặc nhảy qua giải pháp tối ưu. Quá Nhỏ: Dẫn đến tốc độ hội tụ chậm và có thể bị mắc kẹt trong các điểm địa phương. Điều Chỉnh Động: Một số kỹ thuật như learning rate schedules hoặc adaptive learning rates (như Adagrad) có thể được sử dụng để điều chỉnh tỷ lệ học theo thời gian hoặc theo từng trọng số.

SGD là một công cụ mạnh mẽ khi được điều chỉnh đúng cách và thường được kết hợp với các kỹ thuật như Momentum để cải thiện hiệu suất.

## 3. Adam (Adaptive Moment Estimation):

Adam (Adaptive Moment Estimation) là một phương pháp tối ưu hóa được thiết kế để cải thiện quá trình huấn luyện mô hình học máy. Nó kết hợp cả yếu tố động (momentum) và tỷ lệ học thích ứng để tối ưu hóa hiệu suất. Dưới đây là giới thiệu chi tiết về cách Adam hoạt động:

- Tính Gradient: Tính gradient của hàm mất mát đối với từng trọng số của mô hình, tương tự như SGD.
- Tính First-order Moment (Momentum): Tính first-order moment, tương tự như trong phương pháp Momentum.

- Tính Second-order Moment (Độ Lớn của Gradient): Tính second-order moment, là trung bình bình phương của độ lớn của gradient.
- Hiệu Chỉnh Moment và Độ Lớn: Để khắc phục sự chệch đầu tiên khi khởi tạo moment.
- Cập Nhật Trọng Số: Cập nhật trọng số theo hướng ngược của gradient được điều chỉnh bằng moment và độ lớn

Adam là một trong những phương pháp Optimizer phổ biến và được ưa chuộng trong cộng đồng học máy, và nó thường được sử dụng mặc định cho nhiều mô hình do sự linh hoạt và hiệu quả của nó.

#### **4. RMSprop (Root Mean Square Propagation):**

RMSprop (Root Mean Square Propagation) là một phương pháp tối ưu hóa được thiết kế để cải thiện quá trình huấn luyện mô hình học máy bằng cách điều chỉnh tỷ lệ học dựa trên độ lớn của gradient. RMSprop giúp ứng phó tốt với biến động của độ lớn gradient và ngăn chặn các vấn đề liên quan đến tỷ lệ học cố định trong Stochastic Gradient Descent (SGD). Dưới đây là giới thiệu chi tiết về cách RMSprop hoạt động:

- Tính Gradient: Tính gradient của hàm mất mát đối với từng trọng số của mô hình, giống như trong SGD.
- Tính Trung Bình Bình Phương của Độ Lớn Gradient: Tính trung bình bình phương của độ lớn gradient theo từng trọng số
- Cập Nhật Trọng Số: Cập nhật trọng số của mô hình theo hướng ngược của gradient được điều chỉnh bằng tỷ lệ học thích ứng
- Ưu Điểm về Tỷ Lệ Học Thích Ứng: RMSprop thích ứng tỷ lệ học theo từng trọng số. Điều này giúp giảm bớt vấn đề về việc chọn tỷ lệ học cố định trong SGD.
- Hiệu Chỉnh Moment: Tương tự như trong Adam, RMSprop có thể thực hiện việc hiệu chỉnh moment để khắc phục sự chệch đầu tiên khi khởi tạo moment.

RMSprop là một trong những phương pháp tối ưu hóa phổ biến và thường được sử dụng trong thực tế để cải thiện quá trình huấn luyện mô hình học máy.

#### **5. Adagrad (Adaptive Gradient Algorithm):**

Adagrad (Adaptive Gradient Algorithm) là một phương pháp tối ưu hóa được thiết kế để tự động điều chỉnh tỷ lệ học cho từng trọng số của mô hình dựa trên lịch sử của gradient. Adagrad có khả năng tự thích ứng với độ quan trọng của từng tham số, giúp nó trở nên hiệu quả đối với các tham số có độ nhạy khác nhau. Dưới đây là giới thiệu chi tiết về cách Adagrad hoạt động:

- Tính Gradient: Tính gradient của hàm mất mát đối với từng trọng số của mô hình, giống như trong SGD.
- Tính Tổng Bình Phương của Gradient: Tính tổng bình phương của gradient theo từng trọng số.

- Cập Nhật Trọng Số: Cập nhật trọng số của mô hình theo hướng ngược của gradient được điều chỉnh bằng tỷ lệ học thích ứng
- Ưu Điểm về Tự Động Điều Chỉnh Tỷ Lệ Học: Adagrad tự động điều chỉnh tỷ lệ học theo từng trọng số dựa trên lịch sử của gradient. Nó giúp giảm bớt vấn đề về việc chọn tỷ lệ học cố định trong SGD.

Adagrad thường được sử dụng trong các bài toán mà các tham số của mô hình có độ nhạy khác nhau và yêu cầu một sự điều chỉnh tỷ lệ học linh hoạt. Tuy nhiên, để tránh các vấn đề về tích tổng bình phương quá nhanh, các phương pháp như RMSprop và Adam đã được phát triển và thường được ưa chuộng hơn trong nhiều trường hợp.

## 6. Momentum:

Momentum là một phương pháp tối ưu hóa trong quá trình huấn luyện mô hình học máy, được thiết kế để giúp quá trình cập nhật trọng số của mô hình trở nên ổn định hơn và giảm hiện tượng dao động (oscillation) trong quá trình học.

Dưới đây là giới thiệu chi tiết về cách Momentum hoạt động:

- Tính Gradient: Tính gradient của hàm mất mát đối với từng trọng số của mô hình, giống như trong SGD.
- Tính Động Moment (Momentum): Sử dụng động moment để giữ lại một phần của động lượng từ các bước cập nhật trước đó
- Cập Nhật Trọng Số: Cập nhật trọng số của mô hình theo hướng ngược của gradient được điều chỉnh bằng động

Momentum thường được sử dụng như một thành phần quan trọng của các phương pháp tối ưu hóa phức tạp hơn như Adam và Nesterov Accelerated Gradient (NAG). Sự kết hợp này giúp tối ưu hóa hiệu suất và đồng thời giảm thiểu nhược điểm của từng phương pháp riêng lẻ.

## 7. Nadam (Nesterov-accelerated Adaptive Moment Estimation):

Nadam (Nesterov-accelerated Adaptive Moment Estimation) là một phương pháp tối ưu hóa là sự kết hợp giữa phương pháp Nesterov Accelerated Gradient (NAG) và Adaptive Moment Estimation (Adam). Nadam thường được sử dụng trong quá trình huấn luyện mô hình học máy để cải thiện tốc độ hội tụ và ổn định quá trình cập nhật trọng số. Dưới đây là giới thiệu chi tiết về cách Nadam hoạt động

- Tính Gradient: Tính gradient của hàm mất mát đối với từng trọng số của mô hình, giống như trong SGD.
- Tính Moment (Nesterov Acceleration): Sử dụng động moment theo phương pháp Nesterov để điều chỉnh gradient trước khi tính toán độ lớn của gradient.
- Tính Second-order Moment (Adam-like): Tính second-order moment giống như trong Adam.

- Hiệu Chỉnh Moment và Độ Lớn: Hiệu chỉnh động moment và second-order moment.
- Cập Nhật Trọng Số: Cập nhật trọng số của mô hình theo hướng ngược của gradient được điều chỉnh bằng động moment và độ lớn.

Nadam là một trong những phương pháp tối ưu hóa hiệu quả và thường được sử dụng trong các quá trình huấn luyện mô hình học máy đối với các tập dữ liệu và kiến trúc mô hình đa dạng.

## 8. So sánh ưu điểm và nhược điểm của các phương pháp

### Optimizer trên:

Phương pháp Optimizer	Ưu điểm	Nhược điểm	Thích hợp cho
Stochastic Gradient Descent (SGD)	- Đơn giản, dễ triển khai.	- Có thể hội tụ chậm trên các bề mặt hàm phức tạp	- Dữ liệu lớn, mô hình đơn giản
Adam (Adaptive Moment Estimation)	- Hiệu quả cho nhiều loại mô hình	- Cần điều chỉnh các siêu tham số	- Dữ liệu biến động, mô hình đa dạng
RMSprop (Root Mean Square Propagation)	- Hiệu quả ứng phó với biến động của độ lớn gradient	- Cần điều chỉnh các siêu tham số	- Mô hình dữ liệu biến động
Adagrad (Adaptive Gradient Algorithm)	- Hiệu quả cho các tham số có độ quan trọng khác nhau	- Tích tụ độ lớn của gradient có thể làm giảm tỉ lệ học	- Mô hình với các tham số quang trọng không đồng đều
Momentum	- Giúp tránh bẫy địa phương, tăng tốc hội tụ	- Cần điều chỉnh tỉ lệ học và độ lớn của moment	- Mô hình với bề mặt hàm phức tạp, dữ liệu biến động
Nadam (Nesterov-accelerated Adaptive Moment Estimation)	- Kết hợp động moment và tỷ lệ học thích ứng	- Cần điều chỉnh các siêu tham số	- Mô hình đa dạng, dữ liệu biến động, bề mặt hàm phức tạp

## II. Continual learning và Test Production:

### 1. Continual Learning:

Continual Learning (hoặc còn gọi là Lifelong Learning) là một lĩnh vực trong machine learning (học máy) tập trung vào việc phát triển các mô hình có khả năng liên tục học từ dữ liệu mới mà không quên kiến thức đã học từ dữ liệu cũ. Trong quá trình triển khai giải pháp học máy cho một bài toán nào đó, đặc biệt là trong môi trường thay đổi và thêm mới dữ liệu liên tục, Continual Learning

trở thành một yếu tố quan trọng. Dưới đây là một số điểm quan trọng về Continual Learning:

- Memory và Catastrophic Forgetting: Vấn đề quan trọng nhất trong Continual Learning là nguy cơ quên thông tin từ dữ liệu cũ khi học từ dữ liệu mới, được gọi là "catastrophic forgetting." Mô hình cần có khả năng giữ lại và tái sử dụng kiến thức trước đó.
- Regularization và Replay: Các kỹ thuật regularization và replay được sử dụng để giảm thiểu nguy cơ quên thông tin. Replay là quá trình sử dụng dữ liệu cũ để làm mới mô hình, trong khi regularization áp đặt các ràng buộc để bảo vệ kiến thức cũ.
- Incremental và Online Learning: Continual Learning thường thực hiện thông qua các phương pháp incremental hoặc online learning, cho phép mô hình được cập nhật khi có dữ liệu mới mà không cần huấn luyện lại toàn bộ mô hình từ đầu.
- Transfer Learning và Meta-Learning: Transfer learning và meta-learning cũng được tích hợp trong Continual Learning để chia sẻ kiến thức từ nhiệm vụ cũ sang nhiệm vụ mới và nhanh chóng thích ứng với dữ liệu mới.

## 2. Test Production:

Test Production là một phần quan trọng của quá trình triển khai giải pháp học máy sau khi mô hình đã được huấn luyện và kiểm định trên tập dữ liệu đào tạo và kiểm thử. Test Production bao gồm các bước sau:

- Testing on New Data: Kiểm thử mô hình trên dữ liệu mới, không được sử dụng trong quá trình huấn luyện, để đánh giá hiệu suất thực tế của mô hình.
- Performance Metrics: Đánh giá mô hình bằng các chỉ số hiệu suất phù hợp với bài toán cụ thể, chẳng hạn như accuracy, precision, recall, F1-score, hay các metric khác tùy thuộc vào mục tiêu của bài toán.
- Monitoring và Maintenance: Thiết lập các hệ thống giám sát để theo dõi hiệu suất của mô hình trong môi trường thực tế. Nếu có sự chệch lệch hoặc suy giảm hiệu suất, cần có quy trình bảo trì để cập nhật mô hình hoặc thực hiện điều chỉnh.
- Scalability và Deployment: Đảm bảo tính mở rộng và khả năng triển khai của mô hình để xử lý số lượng lớn dữ liệu và đáp ứng nhanh chóng với các yêu cầu sản xuất.
- Feedback Loop: Xây dựng một chu kỳ phản hồi để cập nhật mô hình khi có sự thay đổi trong dữ liệu hoặc yêu cầu kỹ thuật.
- Versioning: Quản lý phiên bản của mô hình để dễ dàng theo dõi các thay đổi, thuận tiện cho việc quay trở lại phiên bản trước đó nếu cần thiết.