



HCMUS

Viet Nam National University
Ho Chi Minh City
University of Science



Khoa Toán - Tin học
Fac. of Math. & Computer Science

CHƯƠNG 4

Neural Networks

Lưu Giang Nam

Bộ môn Ứng dụng Tin học
Khoa Toán - Tin học
Trường Đại học KHTN, ĐHQG TPHCM

09/2025

Bối cảnh và vấn đề

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Trong Chương 3–4: xem xét các mô hình hồi quy và phân loại dựa trên tổ hợp tuyến tính của các hàm cơ sở cố định.
- Ưu điểm: giải tích tốt, tính toán thuận lợi.
- Hạn chế: chịu ảnh hưởng của lời nguyên chiều.
- Mạng nơ-ron truyền thẳng (Feed-forward NN)
 - Cố định số lượng hàm cơ sở, nhưng cho phép chúng thích nghi thông qua các tham số huấn luyện.
 - Mô hình thành công nhất: Feed-forward neural network hay Multilayer Perceptron (MLP).
 - Thực chất MLP là chuỗi các mô hình logistic liên tục, không phải các perceptron rời rạc.
 - Ưu điểm: mô hình gọn hơn, tính toán nhanh hơn SVM khi có cùng năng lực tổng quát hóa.



Phần 1

Feed-forward Network

Mô hình tuyến tính cơ sở

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Bắt đầu từ mô hình tuyến tính hồi quy hoặc phân loại:

$$y(x, w) = f \left(\sum_{j=1}^M w_j \phi_j(x) \right) \quad (1)$$

- Trong đó:
 - $f(\cdot)$ là hàm kích hoạt (activation function).
 - $\phi_j(x)$ là các hàm cơ sở cố định
- Mục tiêu: mở rộng mô hình bằng cách cho phép $\phi_j(x)$ phụ thuộc tham số, huấn luyện cùng w_j
- Neural Network: mỗi $\phi_j(x)$ là một hàm phi tuyến của tổ hợp tuyến tính các đầu vào (inputs) với trọng số (weights) và độ lệch (bias).

Hàm kích hoạt

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Ở mỗi tầng ẩn (hidden unit) ta định nghĩa a_j :

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}, \quad j = 1, \dots, M \quad (2)$$

- Từ dạng linear ta chuyển qua dạng phi tuyến qua thông qua hàm kích hoạt (khả vi, phi tuyến) ta được hàm:

$$z_j = h(a_j) \quad (3)$$

- Các hàm kích hoạt phổ biến:
 - Logistic sigmoid: $\sigma(a) = \frac{1}{1+\exp(-a)}$
 - Tanh: $h(a) = \tanh(a)$
- Chức năng:
 - Tạo phi tuyến giữa input và output
 - Đảm bảo khả vi

Lớp đầu ra

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Tầng đầu ra (output) với giá trị kích hoạt đầu ra ở tầng thứ 2:

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)}, \quad k = 1, \dots, K$$

- Cuối cùng, các kích hoạt đầu ra được biến đổi bằng cách sử dụng hàm kích hoạt thích hợp để đưa ra một tập hợp các đầu ra y_k .
 - Regression: identity, $y_k = f(a_k) = a_k$
 - Binary classification: sigmoid, $y_k = f(a_k) = \sigma(a_k)$
 - Multiclass: softmax, $y_k = f(a_k) = \text{softmax}(a_k)$
- Lan truyền thuận: tuần tự từ input \rightarrow hidden \rightarrow output

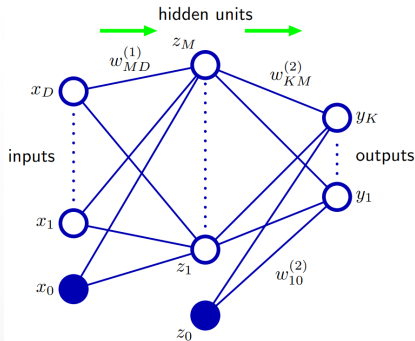
$$y_k(x, w) = f \left(\sum_{j=0}^M w_{kj}^{(2)} h \left(\sum_{i=0}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right) \quad (4)$$

Sơ đồ mạng cho mạng nơ-ron hai lớp

Feed-forward Network

Huấn luyện mô hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent



Hình: Các biến đầu vào, biến ẩn và biến đầu ra được biểu diễn bằng các nút, và các trọng số được biểu diễn bằng các liên kết giữa các nút, trong đó các tham số độ lệch được biểu thị bằng các liên kết đến từ các biến đầu vào và biến ẩn bổ sung x_0 và z_0 . Các mũi tên biểu thị hướng của luồng thông tin qua mạng trong quá trình lan truyền thuận (forward propagation).

Bài tập 1

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

Câu hỏi

Với điều kiện $x_0 = 1$, hãy biểu diễn a_j và $y_k(x, w)$ theo dạng sau:

$$a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i, \quad y_k(x, w) = \sigma \left(\sum_{j=0}^M w_{kj}^{(2)} h \left(\sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right) \quad (5)$$

- Lợi ích:
 - Viết gọn công thức
 - Chuẩn hóa forward propagation

MLP và kiến trúc mạng

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Hai tầng = Multilayer Perceptron (MLP)
- Hidden units liên tục, khả vi \rightarrow hỗ trợ gradient-based training
- Linear hidden units \rightarrow khả năng xấp xỉ hạn chế
- Nonlinear hidden units: \rightarrow xấp xỉ mọi hàm liên tục trên miền compact. Số lượng hidden units quyết định độ chính xác
- Mỗi hidden unit: tính tuyến tính \rightarrow phi tuyến \rightarrow chuyển sang tầng đầu ra.
- Feed-forward: không có cycles, đảm bảo output xác định
- Mạng có thể mở rộng:
 - Thêm layers
 - Skip connections (input \rightarrow output)
 - Sparse connections

Vai trò hidden units

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Hidden units trích xuất đặc trưng phi tuyến
- Hợp tác để xấp xỉ hàm mục tiêu
- Cấu trúc hai tầng + nonlinear hidden units → xấp xỉ nhiều hàm phức tạp
- Thách thức: tìm tập tham số w tối ưu từ dữ liệu huấn luyện
- Giải pháp:
 - Maximum likelihood
 - Bayesian approach
- Mở rộng: Skip-layer connections, sparse connections
- Có thể thêm layers → deep network nhưng luôn tuân thủ feed-forward → output xác định
- Mỗi unit:

$$z_k = h \left(\sum_j w_{kj} z_j \right) \quad (6)$$



Phần 2

Huấn luyện mô hình

Giới thiệu

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Neural network là một lớp hàm phi tuyến tham số từ input vector $x \rightarrow$ output vector y
- Mục tiêu: tìm các tham số w sao cho mạng mô hình hóa tốt dữ liệu huấn luyện
- Phương pháp đơn giản: Cực tiểu hóa tổng bình phương lỗi:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2$$

- Nhưng cách tiếp cận tổng quát hơn là dùng góc nhìn xác suất cho outputs. Làm rõ ràng hơn cho cả việc lựa chọn tính phi tuyến tính của output unit và lựa chọn hàm lỗi.



Probabilistic Interpretation - Regression

- Giả sử một target $t \in \mathbb{R}$ có phân phối Gaussian với

$$p(t|x, w) = \mathcal{N}(t|y(x, w), \beta^{-1}) \quad (7)$$

- β = độ chính xác (precision hay inverse variance) của nhiễu Gauss.
- Đối với phân phối có điều kiện được cho bởi (7), chỉ cần lấy hàm kích hoạt đơn vị đầu ra là đồng nhất thức, vì một mạng như vậy có thể xấp xỉ bất kỳ hàm liên tục nào từ x đến y .
- Cho một tập dữ liệu gồm N quan sát độc lập, phân phối đồng nhất $X = \{x_1, \dots, x_N\}$, cùng với các giá trị mục tiêu tương ứng $t = \{t_1, \dots, t_N\}$ ta có thể xây dựng được likelihood cho toàn bộ dataset:

$$p(t|X, w, \beta) = \prod_{n=1}^N p(t_n|x_n, w, \beta) \quad (8)$$

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent



Bài tập 2:

Câu hỏi

- 1 Chứng minh hàm sai số (error function) được cho bằng hàm đối của log-likelihood là (với w và β là các tham số):

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi) \quad (9)$$

- 2 Sau khi tìm được w_{ML} , giá trị của β có thể được tìm thấy bằng cách tương tự.

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{ML}) - t_n\}^2 \quad (10)$$

Nhiên biến mục tiêu

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Giả sử K biến mục tiêu (independent conditionally) trên x và w , khi đó ta xét:

$$p(t|x, w) = \mathcal{N}(t|y(x, w), \beta^{-1}I) \quad (11)$$

- Khi đó để tối đại likelihood ta sẽ tối thiểu $E(w)$:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \{y_k(x_n, w) - t_{nk}\}^2 \quad (12)$$

với

$$\frac{1}{\beta_{ML}} = \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \{y_k(x_n, w_{ML}) - t_{nk}\}^2 \quad (13)$$

Với hàm đơn nhất

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Nếu Output activation = identity, tức là $y_k = a_k$
- Ta sẽ có:

$$\frac{\partial E}{\partial a_k} = y_k - t_k \quad (14)$$

- Tương thích tự nhiên giữa output activation và error function
- Max likelihood $w \leftrightarrow$ minimize $E(w)$:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 \quad (15)$$

- Nonlinear network $\rightarrow E(w)$ nonconvex \rightarrow local minima

Với phân loại nhị phân

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Với hàm mục tiêu đơn biến t : lớp C_1 gán $t = 1$, lớp C_2 gán $t = 0$.
- Sử dụng kích hoạt cho hàm đầu ra là hàm logistic sigmoid:

$$y = \sigma(a) = \frac{1}{1 + \exp(-a)} \quad (16)$$

- Khi đó ta có $0 \leq y(x, w) \leq 1$ và $y(x, w) = p(C_1|x)$, $p(C_2|x) = 1 - y(x, w)$.
- Khi đó phân phối của điều kiện của biến mục tiêu sẽ tuân theo phân phối Bernoulli

$$p(t|x, w) = y(x, w)^t \{1 - y(x, w)\}^{1-t} \quad (17)$$

Cross-entropy

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Khi đó, sử dụng “negative log-likelihood” ta sẽ được một hàm lỗi được gọi là “cross-entropy”:

$$E(w) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (18)$$

- Không có noise precision β
- Ưu điểm: huấn luyện nhanh hơn, tổng quát hóa tốt hơn.

Cross-entropy

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

Đối với bài toán đa biến (K biến) ta đơn giản là tích (K) hàm sigmoid lại:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \prod_{k=1}^K y_k^{t_k}(\mathbf{x}, \mathbf{w})(1 - y_k(\mathbf{x}, \mathbf{w}))^{1-t_k} \quad (19)$$

Sử dụng “negative log-likelihood” để chứng minh hàm lỗi $E(\mathbf{w})$ sẽ có dạng:

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K \{t_{nk} \ln y_{nk} + (1 - t_{nk}) \ln(1 - y_{nk})\} \quad (20)$$

với $y_{nk} = y_K(\mathbf{x}_n, \mathbf{w})$



Phần 2

Huấn luyện mô hình

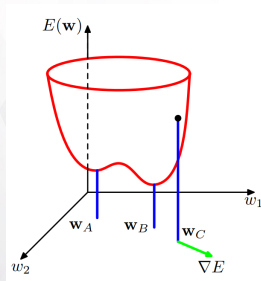
Mục 1: TỐI ƯU THAM SỐ

Giới thiệu

- Mục tiêu: tìm vector trọng số w để **minimize** hàm lỗi $E(w)$.
- Phương pháp được sử dụng ở đây sẽ là cập nhật các trọng số để hội tụ về giá trị cực trị:

$$w \rightarrow w + \delta w, \quad \delta E \simeq \delta w^T \nabla E(w) \quad (21)$$

với vector $\nabla E(w)$ chỉ hướng tăng nhanh nhất của $E(w)$.



Hình: Điểm w_A là một điểm cực tiểu cục bộ và w_B là điểm cực tiểu toàn cục. Điểm w_C và vectơ $\nabla E(w)$ ứng với điểm này.

Điều kiện cực trị của hàm lỗi

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số

Xấp xỉ bậc hai cục bộ

Tối ưu gradient descent

- Các điểm mà $\nabla E(w) = 0$ gọi là **stationary points**.
- Phân loại:
 - Trong mạng nơ-ron, $E(w)$ thường **phi tuyến mạnh** \rightarrow có nhiều điểm dừng.
 - $E(w)$ có thể có nhiều cực tiểu không tương đương:

$$E(w_{\text{local}}) > E(w_{\text{global}}) \quad (22)$$

- Trong thực tế, ta không cần tìm chính xác global minimum \rightarrow tìm nghiệm tốt (local minimum “đủ tốt”) có khả năng tổng quát hóa tốt.

Thuật toán tối ưu lặp

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent

- Trong giải tích, việc tìm nghiệm chính xác của phương trình $\nabla E(w) = 0$ là không thể \rightarrow dùng **thuật toán lặp**.
- Công thức cập nhật tổng quát:

$$w^{(\tau+1)} = w^{(\tau)} + \Delta w^{(\tau)} \quad (23)$$

- Sau mỗi bước, tính lại gradient tại $w^{(\tau+1)}$.
- Các thuật toán khác nhau ở cách chọn $\Delta w^{(\tau)}$.

Phần 2

Huấn luyện mô hình

Mục 2: XẤP XỈ BẬC HAI CỤC BỘ

Local Quadratic Approximation

- Gần một điểm w_b , khai triển Taylor bậc hai:

$$E(w) \simeq E(w_b) + (w - w_b)^T b + \frac{1}{2} (w - w_b)^T H (w - w_b) \quad (24)$$

với:

$$b = \nabla E|_{w=w_b}, \quad H = \nabla \nabla E|_{w=w_b} = \frac{\partial E}{\partial w_i \partial w_j} \Big|_{w=w_b}$$

- Đây là xấp xỉ bậc hai của $E(w)$ quanh w_b . Khi đó gradient tương ứng:

$$\nabla E \simeq b + H(w - w_b) \quad (25)$$

- Nếu w_b là điểm cực tiểu $\rightarrow b = 0$:

$$E(w) = E(w_b) + \frac{1}{2} (w - w_b)^T H (w - w_b) \quad (26)$$

Giá trị riêng và hình dạng bề mặt lồi

- Giải phương trình riêng của ma trận Hessian H :

$$Hu_i = \lambda_i u_i, \quad u_i^T u_j = \delta_{ij} \quad (27)$$

- Bây giờ chúng ta viết lại $(w - w_b)$ thành một tổ hợp tuyến tính của các vectơ riêng dưới dạng: $w - w_b = \sum_i \alpha_i u_i$.

Câu hỏi

Thay vào chứng minh $E(w)$ có dạng:

$$E(w) = E(w_b) + \frac{1}{2} \sum_i \lambda_i \alpha_i^2 \quad (28)$$

Câu hỏi

Khi đó các $\lambda_i > 0$ ta sẽ có ma trận Hessian H là **positive definite**. Chứng minh w_b là cực tiểu.

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent



Phần 2

Huấn luyện mô hình

Mục 3: TỐI ƯU GRADIENT DESCENT

Batch Version

- Quy tắc cập nhật cơ bản:

$$w^{(\tau+1)} = w^{(\tau)} - \eta \nabla E(w^{(\tau)})$$

- $\eta > 0$: learning rate (bước nhảy).
- Mỗi bước cần duyệt qua toàn bộ training set.
- Gọi là batch gradient descent.
- Nhược điểm của Gradient Descent: Dễ dao động khi hàm lỗi có độ cong khác nhau. Dễ mắc kẹt trong local minima (hoặc plateau).
- Cần chọn η phù hợp:
 - η quá nhỏ \rightarrow hội tụ chậm.
 - η quá lớn \rightarrow dao động, không hội tụ.
- Các phương pháp hiệu quả hơn: **Conjugate Gradient**, **Quasi-Newton**.

Feed-forward
Network

Huấn luyện mô
hình

Tối ưu tham số
Xấp xỉ bậc hai cục bộ
Tối ưu gradient descent



On-line (Stochastic) Gradient Descent¹

- Hàm lỗi toàn bộ:

$$E(w) = \sum_{n=1}^N E_n(w) \quad (29)$$

- Cập nhật theo từng mẫu:

$$w^{(\tau+1)} = w^{(\tau)} - \eta \nabla E_n(w^{(\tau)}) \quad (30)$$

- Mỗi bước chỉ dùng một điểm dữ liệu.
- Còn gọi là **stochastic gradient descent (SGD)**.
- Ưu điểm của On-line Learning
 - 1 Xử lý tập dữ liệu lớn hiệu quả hơn.
 - 2 Có khả năng thoát khỏi local minima do nhiễu trong cập nhật.
 - 3 Có thể cập nhật liên tục khi có dữ liệu mới.

¹Le Cun, et al. (1989). Backpropagation applied to handwritten zip code recognition. Neural Computation 1(4), 541–551.

The background features a complex, abstract design. It consists of several large, overlapping circles in various shades of gray and white. Overlaid on these circles are patterns of small, dark dots, creating a halftone or dot-matrix effect. The dots are arranged in a grid-like fashion, with some areas being denser than others, creating a sense of depth and texture. The overall composition is clean and modern, with a focus on geometric shapes and patterns.

Good luck!