

	VIETTEL AI RACE	Public 104
	GIỚI THIỆU VỀ HIDDEN MARKOV MODEL	Lần ban hành: 1

## 1. Khái niệm cơ bản về Hidden Markov Model (HMM)

Hidden Markov Model (HMM) là một mô hình thống kê được sử dụng để phân tích các chuỗi dữ liệu có tính chất tuần tự, trong đó trạng thái thực của hệ thống (trạng thái ẩn) không thể quan sát trực tiếp, nhưng có thể suy ra thông qua các quan sát (observations). HMM kết hợp hai quá trình ngẫu nhiên:

- Một quá trình Markov ẩn, mô tả sự chuyển đổi giữa các trạng thái ẩn.
- Một quá trình phát xạ, liên kết mỗi trạng thái ẩn với một tập các quan sát theo một phân phối xác suất.

HMM thường được biểu diễn thông qua các thành phần sau:

- **Tập trạng thái ẩn (Hidden States):** Đại diện cho các trạng thái không quan sát được của hệ thống.
- **Ma trận xác suất chuyển trạng thái (State Transition Matrix):** Xác định xác suất chuyển từ một trạng thái ẩn này sang trạng thái ẩn khác.
- **Ma trận xác suất phát xạ (Emission Probability Matrix):** Mô tả xác suất của một quan sát cụ thể dựa trên trạng thái hiện tại.
- **Phân phối xác suất ban đầu (Initial State Distribution):** Xác định trạng thái khởi đầu của hệ thống.

### 1.1. Sự khác biệt giữa Markov Chain và HMM

Markov Chain là một mô hình toán học đơn giản hơn HMM, trong đó:

- Trạng thái của Markov Chain là có thể quan sát trực tiếp.
- Xác suất chuyển trạng thái chỉ phụ thuộc vào trạng thái hiện tại, không quan tâm đến các trạng thái trước đó.

Ngược lại, HMM phức tạp hơn:

- Trạng thái ẩn của HMM không thể quan sát trực tiếp, mà chỉ có thể suy đoán thông qua các quan sát.
- HMM bổ sung thêm quá trình phát xạ, liên kết các trạng thái ẩn với dữ liệu quan sát.

Ví dụ minh họa: Trong Markov Chain, nếu ta đang xem một chuỗi các điều kiện thời tiết (nắng, mưa), bạn có thể quan sát trực tiếp điều kiện thời tiết tại từng thời điểm. Trong HMM, các điều kiện thời tiết có thể được ẩn (không trực tiếp quan sát được), nhưng ta có thể suy luận từ các quan sát như mức độ ẩm, nhiệt độ, hoặc áp suất không khí.

	<b>VIETTEL AI RACE</b>	Public 104
	<b>GIỚI THIỆU VỀ HIDDEN MARKOV MODEL</b>	Lần ban hành: 1

## 1.2. Vai trò và ứng dụng của HMM trong các bài toán thực tiễn

HMM đóng vai trò quan trọng trong nhiều lĩnh vực nghiên cứu và ứng dụng, đặc biệt là trong xử lý chuỗi dữ liệu. Một số ứng dụng điển hình của HMM bao gồm:

### 2.1.1. Xử lý ngôn ngữ tự nhiên (Natural Language Processing - NLP):

- Gắn thẻ từ loại (POS Tagging): Dự đoán nhãn ngữ pháp (danh từ, động từ,...) của các từ trong câu.
- Nhận dạng thực thể (Named Entity Recognition): Xác định tên riêng, địa danh, hoặc tổ chức trong văn bản.

### 2.1.2. Nhận dạng giọng nói (Speech Recognition):

- Mô hình hóa các chuỗi âm thanh để chuyển đổi thành văn bản.

### 2.1.3. Phân tích sinh học (Bioinformatics):

- Dự đoán cấu trúc protein từ chuỗi axit amin.
- Phân tích trình tự DNA để xác định gen.

### 2.1.4. Phát hiện bất thường (Anomaly Detection):

- Dự đoán lỗi trong hệ thống máy tính hoặc mạng lưới.
- Phát hiện gian lận trong các giao dịch tài chính.

### 2.1.5. Ứng dụng trong thời gian thực:

- Phân tích dữ liệu cảm biến trong hệ thống IoT (Internet of Things).
- Dự đoán trạng thái hoạt động trong các hệ thống điều khiển tự động.

Nhờ khả năng kết hợp tính xác suất và dự đoán trạng thái, HMM trở thành một công cụ mạnh mẽ trong việc mô hình hóa các quá trình phức tạp mà các trạng thái ẩn không thể quan sát trực tiếp.

## 2. Cấu trúc cơ bản của Hidden Markov Model

### 2.1. Mô hình Markov và trạng thái ẩn

Hidden Markov Model (HMM) là sự mở rộng của mô hình Markov truyền thống, trong đó trạng thái của hệ thống không thể quan sát trực tiếp, mà chỉ có

	VIETTEL AI RACE	Public 104
	GIỚI THIỆU VỀ HIDDEN MARKOV MODEL	Lần ban hành: 1

thể được suy luận từ các quan sát (emissions). Một hệ thống được mô tả bởi HMM có các trạng thái ẩn liên kết với một tập hợp các quan sát cụ thể thông qua xác suất phát xạ.

Trong HMM, hai quá trình ngẫu nhiên được kết hợp:

- **Quá trình Markov ẩn:** Mô tả sự chuyển đổi giữa các trạng thái ẩn theo xác suất.
- **Quá trình phát xạ:** Liên kết mỗi trạng thái ẩn với các quan sát thông qua phân phối xác suất phát xạ.

HMM thường được biểu diễn dưới dạng một đồ thị có hướng, trong đó các nút là trạng thái và các cạnh thể hiện xác suất chuyển đổi giữa các trạng thái.

## 2.2. Các thành phần chính của HMM

Một HMM được định nghĩa bởi bốn thành phần chính:

### 2.2.1. Tập trạng thái (Hidden States)

Tập trạng thái ẩn của HMM được ký hiệu là  $S=\{S_1, S_2, \dots, S_N\}$ , trong đó:

- $S_i$ : Trạng thái ẩn thứ iii.
- $N$ : Số lượng trạng thái ẩn.

Tại mỗi thời điểm, hệ thống sẽ nằm ở một trong các trạng thái  $S_i$ , nhưng trạng thái này không thể quan sát trực tiếp mà chỉ có thể suy ra từ các quan sát.

Ví dụ: Trong bài toán nhận dạng giọng nói, các trạng thái ẩn có thể là các âm vị (phonemes) mà người nói đang phát âm.

### 2.2.2. Ma trận chuyển trạng thái (State Transition Matrix)

Ma trận chuyển trạng thái, ký hiệu là  $A=[a_{ij}]$ , là một ma trận vuông kích thước  $N \times N$ , trong đó:

- $a_{ij}=P(S_j|S_i)$ : Xác suất chuyển từ trạng thái  $S_i$  sang trạng thái  $S_j$ .
- $\sum_{j=1}^N a_{ij} = 1$ : Tổng các xác suất từ một trạng thái phải bằng 1.

Ma trận A biểu diễn các mối quan hệ giữa các trạng thái ẩn trong mô hình.

	<b>VIETTEL AI RACE</b>	Public 104
	<b>GIỚI THIỆU VỀ HIDDEN MARKOV MODEL</b>	Lần ban hành: 1

Ví dụ: Trong một chuỗi thời tiết, xác suất chuyển từ trạng thái "nắng" sang "mưa" là một phần của ma trận chuyển trạng thái.

### 2.2.3. Ma trận xác suất phát xạ (Emission Probability Matrix)

Ma trận xác suất phát xạ, ký hiệu là  $B = [b_j(k)]$ , là một ma trận kích thước  $N \times M$ , trong đó:

- $b_j(k) = P(O_k | S_j)$ : Xác suất quan sát  $O_k$  xảy ra khi hệ thống ở trạng thái  $S_j$ .
- $O = \{O_1, O_2, \dots, O_M\}$ : Tập các quan sát có thể xảy ra, với  $M$  là số lượng quan sát.

Ma trận  $B$  mô tả mối quan hệ giữa trạng thái ẩn và quan sát.

Ví dụ: Trong bài toán nhận dạng giọng nói, các quan sát có thể là các đặc trưng âm thanh (spectral features) được trích xuất từ tín hiệu âm thanh.

### 2.2.4. Phân phối xác suất ban đầu (Initial State Distribution)

Phân phối xác suất ban đầu, ký hiệu là  $\pi = \{\pi_1, \pi_2, \dots, \pi_N\}$ , trong đó:

- $\pi_i = P(S_i)$ : Xác suất hệ thống bắt đầu ở trạng thái  $S_i$ .
- $\sum_{i=1}^N \pi_i = 1$ : Tổng xác suất của tất cả các trạng thái ban đầu phải bằng 1.

Phân phối  $\pi$  cung cấp thông tin về trạng thái khởi đầu của hệ thống trước khi các quan sát được thực hiện.

## 2.3. Công thức tổng quát của HMM

Một HMM được định nghĩa bởi các tham số  $\lambda = (A, B, \pi)$ , trong đó:

- $A = [a_{ij}]$ : Ma trận chuyển trạng thái.
- $B = [b_j(k)]$ : Ma trận xác suất phát xạ.
- $\pi = \{\pi_i\}$ : Phân phối xác suất ban đầu.

Cho một chuỗi quan sát  $O = \{O_1, O_2, \dots, O_T\}$  với chiều dài  $T$ , xác suất của chuỗi quan sát được tính theo công thức:

$$P(O | \lambda) = \sum_Q P(O, Q | \lambda) = \sum_Q P(Q | \lambda) \cdot P(O | Q, \lambda)$$

Trong đó:

	VIETTEL AI RACE	Public 104
	GIỚI THIỆU VỀ HIDDEN MARKOV MODEL	Lần ban hành: 1

- $Q = \{S_{q1}, S_{q2}, \dots, S_{qT}\}$ : Một chuỗi trạng thái ẩn.
- Tổng  $\Sigma_Q$  được tính trên tất cả các chuỗi trạng thái có thể xảy ra.

Công thức này cho phép ta tính xác suất quan sát của một chuỗi và xác định chuỗi trạng thái ẩn tối ưu.

### 3. Ba bài toán cơ bản của Hidden Markov Model (HMM)

Hidden Markov Model (HMM) được sử dụng để giải quyết ba bài toán cơ bản trong các ứng dụng thực tế. Các bài toán này là trung tâm của việc áp dụng HMM vào việc phân tích dữ liệu tuần tự. Dưới đây là chi tiết từng bài toán.

#### 3.1. Bài toán 1: Đánh giá (Evaluation)

##### Mục tiêu:

Tính xác suất của một chuỗi quan sát  $O = \{O_1, O_2, \dots, O_T\}$  đã cho, dựa trên mô hình HMM  $\lambda = (A, B, \pi)$ .

##### Ý nghĩa:

Bài toán này giúp đánh giá mức độ phù hợp của một chuỗi quan sát với một mô hình HMM cụ thể. Đây là bước cần thiết để so sánh và lựa chọn mô hình tốt nhất từ các mô hình cạnh tranh.

##### Công thức:

Xác suất của chuỗi quan sát  $P(O | \lambda)$  được tính bằng cách tổng hợp xác suất trên tất cả các chuỗi trạng thái ẩn  $Q = \{q_1, q_2, \dots, q_T\}$ :

$$P(O | \lambda) = \sum_Q P(O, Q | \lambda)$$

##### Thách thức:

Việc tính toán trực tiếp rất phức tạp, vì số lượng các chuỗi trạng thái  $Q$  tăng theo hàm mũ với chiều dài  $T$  của chuỗi quan sát.

##### Giải pháp:

Sử dụng thuật toán **Forward**:

- Thuật toán này tính toán xác suất một cách hiệu quả bằng cách sử dụng phương pháp đệ quy.
- Độ phức tạp được giảm từ  $O(N^T)$  xuống  $O(N^2T)$ , trong đó  $N$  là số trạng thái ẩn.

	VIETTEL AI RACE	Public 104
	GIỚI THIỆU VỀ HIDDEN MARKOV MODEL	Lần ban hành: 1

### 3.2. Bài toán 2: Giải mã (Decoding)

#### Mục tiêu:

Tìm chuỗi trạng thái ẩn tối ưu  $Q^*=\{q_1^*, q_2^*, \dots, q_T^*\}$  tương ứng với chuỗi quan sát  $O$ , sao cho:

$$Q^* = \operatorname{argmax}_Q P(Q | O, \lambda)$$

#### Ý nghĩa:

Bài toán này giúp xác định chuỗi trạng thái ẩn khả dĩ nhất, giải thích tốt nhất cho chuỗi quan sát. Đây là một bước quan trọng trong các ứng dụng như nhận dạng giọng nói và phân tích chuỗi sinh học.

#### Thách thức:

Việc tìm kiếm chuỗi trạng thái tối ưu yêu cầu tối ưu hóa toàn cục trên toàn bộ chuỗi thời gian.

#### Giải pháp:

Sử dụng thuật toán **Viterbi**:

- Thuật toán này dựa trên lập trình động, tìm chuỗi trạng thái tối ưu bằng cách lưu trữ các giá trị tối đa tại mỗi bước.
- Độ phức tạp của thuật toán là  $O(N^2T)$ .

### 3.3. Bài toán 3: Học (Learning)

#### Mục tiêu:

Ước lượng các tham số của mô hình  $\lambda=(A, B, \pi)$  từ một tập dữ liệu quan sát  $O=\{O^{(1)}, O^{(2)}, \dots, O^{(K)}\}$ .

#### Ý nghĩa:

Bài toán này giúp xây dựng một mô hình HMM phù hợp từ dữ liệu quan sát, phục vụ cho việc phân tích và dự đoán.

#### Thách thức:

Không thể tối ưu trực tiếp  $P(O|\lambda)$  vì các trạng thái ẩn không được quan sát trực tiếp.

#### Giải pháp:

Sử dụng thuật toán **Baum-Welch** hoặc **Expectation-Maximization (EM)**:

	<b>VIETTEL AI RACE</b>	Public 104
	<b>GIỚI THIỆU VỀ HIDDEN MARKOV MODEL</b>	Lần ban hành: 1

- Thuật toán này lặp lại hai bước:
  1. **E-step (Expectation):** Tính xác suất kỳ vọng cho các trạng thái ẩn dựa trên các tham số hiện tại.
  2. **M-step (Maximization):** Cập nhật các tham số  $A, B, \pi$  để tối đa hóa xác suất quan sát  $P(O|\lambda)$ .
- Thuật toán hội tụ đến một cực đại cục của  $P(O|\lambda)$ .