

# PATTERN RECOGNITION

**Le Dinh Long**

Faculty of Information Technology  
Industrial University of Ho Chi Minh City, Viet Nam

Ngày 26 tháng 2 năm 2025

## CHƯƠNG 3. ƯỚC LƯỢNG DẠNG THAM SỐ

### 3.1 Giới thiệu

Ở chương hai chúng ta đã tìm hiểu cách xây dựng bộ phân lớp tối ưu dựa trên mô hình thống kê nếu xác suất tiên nghiệm  $P(\omega)$  và mật độ xác suất có điều kiện  $p(\mathbf{x}|\omega)$  cho mỗi lớp được xác định. Tuy nhiên, trong phần lớn các ứng dụng nhận dạng mẫu, các đại lượng này không có sẵn. Trong một số trường hợp chúng ta có thông tin rất mơ hồ về nó cùng với tập dữ liệu mẫu hay tập huấn luyện.

Một trong các tiếp cận phổ biến là chúng ta ước lượng các đại lượng  $P(\omega)$  và  $p(\mathbf{x}|\omega)$  từ tập dữ liệu huấn luyện được thu thập, sau đó nó được sử dụng để xây dựng bộ nhận dạng. Trong hai đại lượng trên thì xác suất tiên nghiệm  $P(\omega)$  có thể được ước lượng dễ dàng đối với các bài toán nhận dạng dựa trên mô hình có giám sát. Riêng mật độ xác suất có điều kiện  $p(\mathbf{x}|\omega)$  tương đối phức tạp hơn, đặc biệt trong trường hợp số mẫu dữ liệu nhỏ và số chiều lớn. Nếu dạng của phân phối được biết (ví dụ như phân phối Gauss) khi đó hàm mật độ xác suất được xác định khi các thông số của phân phối (ví dụ như kỳ vọng  $\mu_i$  và ma trận hiệp phương sai  $\Sigma_i$ ) được xác định. Bài toán ước lượng hàm mật độ xác suất đơn giản trở thành bài toán ước lượng các tham số của nó (ví dụ bài toán ước lượng kỳ vọng  $\mu_i$  và ma trận hiệp phương sai  $\Sigma_i$ ).

Đối với bài toán ước lượng tham số, có hai phương pháp phổ biến trong thống kê là ước lượng hợp lý cực đại hay viết tắt là MLE (maximum likelihood estimation) và ước lượng Bayes (BE hay Bayesian estimation). Mặc dù kết quả của hai phương pháp ước lượng tham số này thường gần giống nhau trong nhiều ứng dụng, nhưng cách tiếp cận của nó khá khác nhau. Phương pháp MLE giả thiết rằng các tham số mặc dù không biết nhưng nó được cố định ở đâu đó, và mục tiêu của ước lượng là tìm các tham số này sao cho cực đại xác suất đạt được mẫu thu thập. Trong khi đó tiếp cận BE xem các tham số này là các biến ngẫu nhiên với phân phối cho trước. Cùng với các mẫu thu thập sẽ tạo thành mật độ hậu nghiệm để dần hiệu chỉnh độ chính xác về các tham số cần được ước lượng. Chúng ta sẽ thấy rằng, tiếp cận này sẽ làm hàm mật độ xác suất hậu nghiệm của tham số được nhọn hơn (có độ tập trung cao hơn) xung quanh đại lượng chính xác khi số mẫu tăng. Chi tiết về hai phương pháp ước lượng này được trình bày chi tiết ở các phần sau.

## 3.2 Ước lượng hợp lý cực đại (Maximum likelihood estimation, MLE)

### 3.2.1 Cơ bản về MLE

MLE là một trong các phương pháp ước lượng được quan tâm rất nhiều do khả năng hội tụ của nó khi số mẫu huấn luyện tăng và đơn giản hơn một số phương pháp ước lượng khác như BE.

MLE yêu cầu đầu tiên có tập dữ liệu huấn luyện. Tập mẫu này có thể được tách thành các tập con tương ứng cho từng lớp. Giả sử có  $C$  lớp thì khi đó sẽ có  $C$  tập con được ký hiệu là  $D_1, D_2, \dots, D_C$  với các mẫu trong mỗi tập  $D_i$  được tạo ra một cách độc lập từ mật độ phân phối xác suất  $p(\mathbf{x}|\omega_i)$ ; các mẫu này được xem là các biến ngẫu nhiên phân phối đồng nhất độc lập (independent identically distributed – i.i.d.). Giả sử rằng  $p(\mathbf{x}|\omega_i)$  có dạng tham số hóa, có nghĩa là dạng của phân phối này được biết và được xác định chính xác bởi một vector tham số  $\theta_i$ . Ví dụ nếu như  $p(\mathbf{x}|\omega_i)$  có dạng phân phối Gauss,  $p(\mathbf{x}|\omega_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ , thì tham số  $\theta_i$  bao gồm hai thành phần  $\boldsymbol{\mu}_i$  và  $\boldsymbol{\Sigma}_i$ . Chúng ta có thể sử dụng  $p(\mathbf{x}|\omega_i; \theta_i)$  để thể hiện sự phụ thuộc của  $p(\mathbf{x}|\omega_i)$  với  $\theta_i$ . Chú ý  $\theta_i$  trong trường hợp này không phải là biến ngẫu nhiên. Bài toán ước lượng mật độ xác suất trở thành ước lượng các tham số  $\theta_i$  tương ứng cho từng lớp từ tập dữ liệu huấn luyện.

Để đơn giản, chúng ta giả sử rằng việc ước lượng các tham số  $\theta_i$  cho mỗi lớp  $\omega_i$  với tập dữ liệu  $D_i$  là độc lập nhau; cách thực hiện giống như nhau. Do đó việc ước lượng các tham số này có thể phát biểu thành dạng chung là “ước lượng tham số  $\theta$  của hàm mật độ phân phối  $p(\mathbf{x}; \theta)$  từ tập dữ liệu huấn luyện  $D$ ”.

Giả sử tập dữ liệu huấn luyện  $D$  có  $n$  mẫu được thu nhận từ phân phối đồng nhất độc lập (i.i.d),  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ . Khi đó ta có

$$\begin{aligned} p(D; \theta) &= p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n; \theta) \\ &= \prod_{j=1}^n p(\mathbf{x}_j; \theta) \end{aligned} \quad (3.1)$$

Xác suất này được xem là likelihood của  $\theta$  đối với tập dữ liệu  $D$  và có thể ký hiệu bằng hàm  $L(\cdot)$ , ví dụ  $L(D; \theta)$ , gọi là hàm likelihood.

Mục tiêu của MLE là tìm tập tham số  $\theta$  sao cho cực đại hàm likelihood này, nghĩa là

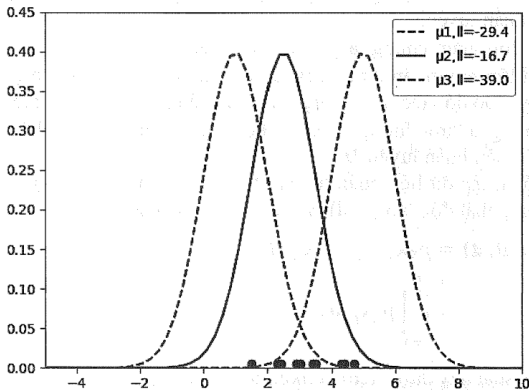
$$\hat{\theta} = \arg \max_{\theta} L(\mathbf{D}; \theta), \quad (3.2)$$

trong đó  $L(\mathbf{D}; \theta)$  là tích của các hàm  $p(\mathbf{x}_j; \theta)$ . Tuy nhiên tích các hàm nhỏ có thể dẫn đến tình trạng không ổn định số học (numerical instability). Một trong cách khắc phục tình trạng này là phát biểu lại bài toán. Chúng ta thấy rằng hàm logarit là hàm đơn điệu tăng, nghĩa là nếu  $\theta$  làm cho likelihood đạt cực đại thì cũng sẽ làm cho log-likelihood đạt cực đại, do đó sử dụng hàm này để phát biểu lại bài toán là một trong những tiếp cận thường được sử dụng. Chúng ta định nghĩa hàm log-likelihood bởi

$$\begin{aligned} l(\theta) &\triangleq \ln L(\mathbf{D}; \theta) \\ &= \sum_{j=1}^n \ln p(\mathbf{x}_j; \theta). \end{aligned} \quad (3.3)$$

Khi đó bài toán trở thành tìm tham số  $\theta$  sao cho cực đại hàm log-likelihood này, nghĩa là

$$\hat{\theta} = \arg \max_{\theta} l(\theta). \quad (3.4)$$



Hình 3.1 Minh họa cho log-likelihood (ll) với tập dữ liệu được cho như trên. Chúng ta thấy rằng với phân phối có  $\mu_2=2.5$  cho giá trị log-likelihood (-16.7) lớn hơn hai trường hợp với phân phối có  $\mu_1=5$  (-29.4) và  $\mu_3=1$  (-39.0).

Tổng quát ta phải tìm lời giải toàn cục cho (3.4). Việc tìm lời giải này có thể được thực hiện theo cách tiếp cận truyền thống dựa trên phép tính vi phân. Giả sử thông số  $\theta$  là vectơ gồm có  $K$  thành phần,  $\theta = (\theta_1, \theta_2, \dots, \theta_K)^T$  và phép tính gradient được định nghĩa bởi

$$\nabla_{\theta} \triangleq \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_K} \end{bmatrix}. \quad (3.5)$$

Khi đó  $\hat{\theta}$  là nghiệm của phương trình

$$\nabla_{\theta} l = 0, \quad (3.6)$$

trong đó  $\nabla_{\theta} l$  được tính bởi

$$\nabla_{\theta} l = \sum_{j=1}^n \nabla_{\theta} \ln p(\mathbf{x}_j; \theta). \quad (3.7)$$



Tóm lại để ước lượng hàm mật độ phân phối dựa trên MLE, chúng ta cần thực hiện bốn bước chính sau:

- Thu thập dữ liệu.
- Chọn mô hình hoặc dạng phân phối, với các tham số có thể được cập nhật.
- Xây dựng hàm mục tiêu dựa trên likelihood, ví dụ như log-likelihood.
- Tối ưu hàm mục tiêu để tìm ra các tham số của mô hình.

### 3.2.2 Trường hợp phân phối chuẩn

Phần này trình bày phương pháp MLE được áp dụng cho trường hợp phân phối cụ thể, đó là phân phối Gauss. Phân phối này có hai tham số là số bình quân (mean)  $\mu$  và ma trận hiệp phương sai  $\Sigma$ . Hàm mật độ phân phối có dạng

$$p(\mathbf{x}; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right), \quad (3.8)$$

trong đó  $d$  là số chiều của dữ liệu. Lấy logarit của  $(\mathbf{x}; \mu, \Sigma)$  ta được

$$\begin{aligned} \ln p(\mathbf{x}; \mu, \Sigma) = & -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] \\ & -\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu). \end{aligned} \quad (3.9)$$

**a) Trường hợp không biết  $\mu$**

Để đơn giản, chúng ta xét trường hợp đầu tiên, chỉ có  $\mu$  là không biết và cần được xác định từ tập dữ liệu  $D$ . Từ (3.7) và (3.9) ta được

$$\nabla_{\mu} l = \sum_{j=1}^n \nabla_{\mu} \ln p(\mathbf{x}_j; \theta) = \sum_{j=1}^n [\Sigma^{-1}(\mathbf{x}_j - \mu)]. \quad (3.10)$$

Thay vào (3.6) ta được

$$\sum_{j=1}^n [\Sigma^{-1}(\mathbf{x}_j - \hat{\mu})] = 0. \quad (3.11)$$

Nhân  $\Sigma$  vào hai vế và thực hiện một số thao tác, ta được kết quả ước lượng của  $\mu$  là

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j. \quad (3.12)$$

Như vậy, MLE của số bình quân của quần thể (population mean) bằng số bình quân mẫu. Nó chính là trung bình cộng của các mẫu huấn luyện.

### b) Trường hợp không biết $\mu$ và $\Sigma$

Trong trường hợp tổng quát hơn, cả mean và ma trận hiệp phương sai đều không biết. Trước tiên chúng ta xét trường hợp một chiều ( $d=1$ ), khi đó (3.8) và (3.9) trở thành

$$p(x; \theta) = \frac{1}{\sqrt{2\pi\theta_2}} \exp\left(-\frac{1}{2\theta_2}(x - \theta_1)^2\right) \quad (3.13)$$

và

$$\ln p(x; \theta) = -\frac{1}{2} \ln[2\pi\theta_2] - \frac{1}{2\theta_2}(x - \theta_1)^2, \quad (3.14)$$

trong đó  $\theta = [\theta_1, \theta_2]^T$ ,  $\theta_1 = \mu$ , và  $\theta_2 = \sigma^2$ . Đạo hàm của nó là

$$\nabla_{\theta} \ln p(x; \theta) = \begin{bmatrix} \frac{1}{\theta_2}(x - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x - \theta_1)^2}{2\theta_2^2} \end{bmatrix}. \quad (3.15)$$

Thay (3.15) vào (3.6) và (3.7) ta có

$$\sum_{j=1}^n \begin{bmatrix} \frac{1}{\theta_2}(x_j - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_j - \theta_1)^2}{2\theta_2^2} \end{bmatrix} = 0. \quad (3.16)$$

Như vậy MLE của  $\theta_1$  và  $\theta_2$  thỏa các phương trình

$$n\hat{\theta}_1 - \sum_{j=1}^n x_j = 0 \quad (3.17)$$

và

$$n\hat{\theta}_2 - \sum_{j=1}^n (x_j - \hat{\theta}_1)^2 = 0 \quad (3.18)$$

Thay  $\hat{\theta}_1 = \hat{\mu}$  và  $\hat{\theta}_2 = \hat{\sigma}^2$  chúng ta được MLE cho  $\mu$  và  $\sigma^2$  lần lượt là

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n x_j \quad (3.19)$$

và

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \hat{\mu})^2. \quad (3.20)$$

Trong trường hợp dữ liệu nhiều chiều ( $d > 1$ ), quá trình cũng được thực hiện tương tự và chúng ta cũng đạt được MLE  $\hat{\mu}$  và  $\hat{\Sigma}$  lần lượt cho  $\mu$  và  $\Sigma$  là

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \quad (3.21)$$

và

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \hat{\mu})(\mathbf{x}_j - \hat{\mu})^T. \quad (3.22)$$

Như vậy, trong trường hợp này ước lượng của vector số bình quân (mean) cũng là số bình quân của tập mẫu (sample mean) và ước lượng của ma trận hiệp phương sai cũng là trung bình cộng của  $n$  ma trận  $(\mathbf{x}_j - \hat{\mu})(\mathbf{x}_j - \hat{\mu})^T$ .

### 3.3 Ước lượng không chệch (unbiased estimation)

Một trong các tiêu chí đánh giá bộ ước lượng là tính không chệch (unbias) của nó. Một bộ ước lượng được gọi là không chệch (unbiased estimator) nếu kỳ vọng của nó bằng giá trị đúng của đại lượng cần ước lượng, nghĩa là nếu  $\hat{\theta}$  là ước lượng của  $\theta$  thì

$$E[\hat{\theta}] = \theta. \quad (3.23)$$

Ước lượng MLE  $\hat{\mu}$  và  $\hat{\Sigma}$  có những đặc điểm về độ không chệch lệch như sau:

- a) Ước lượng MLE  $\hat{\mu}$  của  $\mu$  theo (3.12) và (3.21) là không chệch.
- b) Ước lượng MLE  $\hat{\Sigma}$  của  $\Sigma$  là chệch và ước lượng không chệch của  $\Sigma$  là

$$\hat{\Sigma}_C = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T. \quad (3.24)$$

### 3.4 Phương pháp ước lượng Bayes

Phần này trình bày phương pháp ước lượng phổ biến thứ hai, đó là ước lượng Bayes (BE). Mặc dù kết quả của phương pháp này gần giống như kết quả của MLE. Tuy nhiên tiếp cận của hai phương pháp này khác nhau. Phương pháp MLE xem các đại lượng cần ước lượng là cố định, trong khi phương pháp này xem các đại lượng cần ước lượng là các biến ngẫu nhiên theo một phân phối nào đó. Phương pháp BE sẽ chuyển từ phân phối đó thành phân phối xác suất hậu nghiệm.

Từ công thức Bayes, để tính xác suất hậu nghiệm tương ứng mỗi lớp, ta cần tính các đại lượng xác suất tiên nghiệm  $P(\omega_i)$  và mật độ xác suất có điều kiện  $p(\mathbf{x}|\omega_i)$ . Các đại lượng này được xác định dựa vào thông tin từ tập mẫu  $D$ . Do đó, để nhấn mạnh vai trò của tập mẫu, công thức Bayes có thể được viết thành

$$P(\omega_i|\mathbf{x}, D) = \frac{p(\mathbf{x}|\omega_i, D)P(\omega_i|D)}{\sum_{j=1}^C p(\mathbf{x}|\omega_j, D)P(\omega_j|D)} \quad (3.25)$$

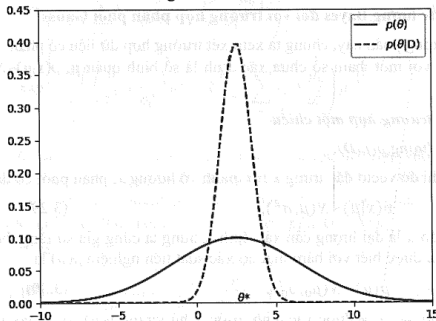


Để đơn giản, chúng ta có thể xem đại lượng  $P(\omega_i|\mathbf{D})$  được xác định hoặc tính dễ dàng từ tập mẫu. Nó có thể được thay bằng  $P(\omega_i)$  và không quan tâm nhiều ở đây. Hơn nữa, tương tự như phương pháp MLE, tập dữ liệu dùng cho huấn luyện có thể được phân thành  $C$  tập con tương ứng với  $C$  lớp. Việc xác định hàm mật độ phân phối cho mỗi lớp không chịu sự ảnh hưởng từ tập dữ liệu của lớp khác, nghĩa là tập con dữ liệu  $\mathbf{D}_j$  sẽ không ảnh hưởng đến việc xác định  $p(\mathbf{x}|\omega_i; \mathbf{D}_i)$  nếu  $i \neq j$ . Do đó công thức (3.25) trở thành

$$P(\omega_i|\mathbf{x}, \mathbf{D}) = \frac{p(\mathbf{x}|\omega_i, \mathbf{D}_i)P(\omega_i)}{\sum_{j=1}^C p(\mathbf{x}|\omega_j, \mathbf{D}_j)P(\omega_j)}. \quad (3.26)$$

Việc ước lượng tham số của các phân phối cho các lớp là tương tự nhau và có thể được thực hiện độc lập nhau. Bài toán trở thành “sử dụng tập dữ liệu huấn luyện  $\mathbf{D}$  từ một phân phối chưa biết (nhưng cố định)  $p(\mathbf{x})$  để ước lượng mật độ phân phối  $p(\mathbf{x}|\mathbf{D})$ ”.

Mặc dù  $p(\mathbf{x})$  là không biết, nhưng dạng của nó thể hiện qua  $p(\mathbf{x}|\theta)$  được xác định và được tham số hóa thông qua  $\theta$ . Bài toán tìm mật độ phân phối trở thành bài toán tìm tham số  $\theta$ . Khác với phương pháp MLE, ở đó tham số của phân phối  $\theta$  là cố định, tham số của phân phối trong trường hợp này được xem là biến ngẫu nhiên và thông tin về nó là xác suất tiên nghiệm với hàm mật độ  $p(\theta)$  được xác định trước. Từ tập dữ liệu huấn luyện thu thập được  $D$  chúng ta xác định xác suất hậu nghiệm  $p(\theta|D)$  được kỳ vọng có đỉnh xấp xỉ và nhọn hơn (sharp) tại giá trị chính xác của  $\theta$  như trong Hình 3.2.



Hình 3.2 Phân phối xác suất tiên nghiệm của tham số  $p(\theta)$ . Khi có tập dữ liệu huấn luyện, nó sẽ làm cho hàm mật độ xác suất được nhọn hơn (với phương sai nhỏ hơn) tại giá trị chính xác của tham số,  $\theta^*$ .

Mục đích cuối cùng của chúng ta trong phương pháp này là xác định  $p(\mathbf{x}|\mathbf{D})$  để xấp xỉ về phân phối đúng  $p(\mathbf{x})$ .  $p(\mathbf{x}|\mathbf{D})$  có thể được tính dựa trên xác suất kết hợp (joint probability)  $p(\mathbf{x},\boldsymbol{\theta}|\mathbf{D})$  như sau

$$p(\mathbf{x}|\mathbf{D}) = \int p(\mathbf{x}, \boldsymbol{\theta}|\mathbf{D}) d\boldsymbol{\theta}. \quad (3.27)$$

Trong đó tích phân được thực hiện trên toàn bộ không gian thông số  $\boldsymbol{\theta}$ . Chúng ta có  $p(\mathbf{x},\boldsymbol{\theta}|\mathbf{D})=p(\mathbf{x}|\boldsymbol{\theta},\mathbf{D})p(\boldsymbol{\theta}|\mathbf{D})$ , hơn nữa  $\mathbf{x}$  và  $\mathbf{D}$  được chọn một cách độc lập nhau, do đó  $p(\mathbf{x}|\boldsymbol{\theta},\mathbf{D})=p(\mathbf{x}|\boldsymbol{\theta})$ . Như vậy (3.27) sẽ trở thành

$$p(\mathbf{x}|\mathbf{D}) = \int p(\mathbf{x}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{D})d\boldsymbol{\theta}. \quad (3.28)$$

Nếu  $p(\boldsymbol{\theta}|\mathbf{D})$  có phương sai rất nhỏ và số bình quân (mean) là  $\hat{\boldsymbol{\theta}}$  thì  $p(\mathbf{x}|\mathbf{D})$  từ (3.28) có thể được xấp xỉ là  $p(\mathbf{x}|\mathbf{D}) \approx p(\mathbf{x}|\hat{\boldsymbol{\theta}})$ . Trong trường hợp tổng quát, giá trị của  $\boldsymbol{\theta}$  không chắc chắn, chúng ta có thể ước lượng  $p(\mathbf{x}|\mathbf{D})$  dựa trên trung bình  $p(\mathbf{x}|\boldsymbol{\theta})$  trên tất cả các giá trị có thể của  $\boldsymbol{\theta}$ . Tích phân (3.28) cũng có thể được tính thông qua một số mô hình, ví dụ như Monte-Carlo.

### 3.5 Ước lượng Bayes đối với trường hợp phân phối Gauss

Trong phần này, chúng ta xem xét trường hợp dữ liệu có phân phối Gauss với một tham số chưa xác định là số bình quân  $\mu$ ,  $p(\mathbf{x}|\mu) \sim N(\mu, \Sigma)$ .

#### a) Xét trường hợp một chiều

\* Ước lượng  $p(\mu|D)$ :

Khi đó vector đặc trưng  $\mathbf{x}$  trở thành vô hướng  $x$ , phân phối có dạng

$$p(x|\mu) \sim N(\mu, \sigma^2), \quad (3.29)$$

trong đó  $\mu$  là đại lượng cần xác định. Chúng ta cũng giả sử rằng thông tin về  $\mu$  được biết với hàm mật độ xác suất tiên nghiệm  $p(\mu)$  là

$$p(\mu) \sim N(\mu_0, \sigma_0^2), \quad (3.30)$$

trong đó  $\mu_0$  và  $\sigma_0$  được xác định trước. Chú ý rằng  $p(\mu)$  có thể có dạng hàm mật độ phân phối khác thay vì Gauss. Với dạng Gauss trên, nó thể hiện giá trị dự đoán tốt nhất của  $\mu$  là  $\mu_0$  và độ lệch chuẩn  $\sigma_0$  thể hiện sự không chắc chắn của dự đoán trên.

Với giá trị  $\mu$  được chọn từ hàm mật độ phân phối  $p(\mu)$ , tập mẫu  $D = \{x_1, x_2, \dots, x_n\}$  được chọn ngẫu nhiên từ hàm mật độ phân phối  $p(x|\mu)$ . Theo công thức Bayes ta có

$$\begin{aligned} p(\mu|D) &= \frac{p(D|\mu)p(\mu)}{\int p(D|\mu)p(\mu)d\mu} \\ &= a \prod_{j=1}^n p(x_j|\mu)p(\mu), \end{aligned} \quad (3.31)$$

trong đó  $a$  là hệ số phụ thuộc tập mẫu  $\mathbf{D}$  nhưng độc lập với  $\mu$ , nó đóng vai trò chuẩn hóa. Thay  $p(x|\mu) \sim N(\mu, \sigma^2)$  và  $p(\mu) \sim N(\mu_0, \sigma_0^2)$  vào (3.31) ta được

$$\begin{aligned}
 p(\mu|\mathbf{D}) &= a \prod_{j=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x_j - \mu}{\sigma}\right)^2\right] \frac{1}{\sigma_0\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right] \\
 &= a_1 \exp\left[-\frac{1}{2}\left(\sum_{j=1}^n \left(\frac{x_j - \mu}{\sigma}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right)\right] \quad (3.32) \\
 &= a_2 \exp\left[-\frac{1}{2}\left(\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2} \sum_{j=1}^n x_j + \frac{\mu_0}{\sigma_0^2}\right)\mu\right)\right],
 \end{aligned}$$

trong đó các thành phần không phụ thuộc  $\mu$  được tích hợp vào các hệ số  $a$ ,  $a_1$ , và  $a_2$ . Rõ ràng  $p(\mu|\mathbf{D})$  có dạng hàm mũ của biểu thức bậc hai theo  $\mu$ , có nghĩa nó cũng là phân phối chuẩn. Nếu ta biểu diễn hàm mật độ phân phối này theo hai tham số  $\mu_n$  và  $\sigma_n$ ,  $p(\mu|\mathbf{D}) \sim N(\mu_n, \sigma_n^2)$ , khi đó nó có dạng

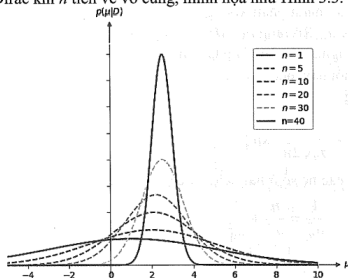
Giải hệ hai phương trình (3.34) và (3.35) ta được

$$\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \bar{x}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \quad (3.37)$$

và

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}. \quad (3.38)$$

Hai phương trình trên cung cấp mối quan hệ giữa xác suất hậu nghiệm của thông số  $p(\mu|\mathbf{D})$  với thông tin xác suất tiên nghiệm  $p(\mu)$  và tập dữ liệu mẫu  $\mathbf{D}$ . Cụ thể nếu không có tập mẫu  $\mathbf{D}$ ,  $n=0$ , thì dự đoán tốt nhất của  $\mu$  sẽ hướng về  $\mu_0$ , với tập dữ liệu  $\mathbf{D}$  có  $n$  mẫu thì ước lượng của  $\mu$  là  $\mu_n$  với độ không chắc chắn được xác định dựa vào  $\sigma_n^2$ . Từ (3.38) dễ dàng thấy rằng  $\sigma_n^2$  đơn điệu giảm khi  $n$  tăng và tiệm cận  $\sigma^2/n$  khi  $n$  tiến về vô cùng. Khi  $n$  càng tăng hàm  $p(\mu|\mathbf{D})$  sẽ càng nhọn hơn và tiệm cận về hàm Dirac khi  $n$  tiến về vô cùng, minh họa như Hình 3.3.



từ xác suất tiên nghiệm là rất cao, nó không làm thay đổi quan điểm giá trị của  $\mu_n$  cho dù chúng ta có thêm tập mẫu. Ngược lại, nếu  $\sigma_0$  rất lớn so với  $\sigma$  ( $\sigma_0 \gg \sigma$ ), có nghĩa độ chắc chắn của ước lượng  $\mu_n = \mu_0$  rất thấp, khi đó ước lượng của  $\mu_n$  hướng về tập dữ liệu thu thập, nghĩa là  $\mu_n \approx \bar{x}_n$ . Trong các trường hợp còn lại, khi số mẫu đủ lớn thì đóng góp của  $\mu_0$  và  $\sigma_0$  sẽ giảm dần và  $\mu_n$  sẽ hội tụ về số bình quân mẫu,  $\bar{x}_n$ .

\* Ước lượng  $p(x|\mathbf{D})$ :

Sau khi ước lượng  $p(\mu|\mathbf{D})$ , công việc còn lại là ước lượng hàm mật độ xác suất có điều kiện  $p(x|\mathbf{D})$ . Từ (3.28) ta có

$$p(x|\mathbf{D}) = \int p(x|\mu) p(\mu|\mathbf{D}) d\mu. \quad (3.39)$$

Thay (3.29) và (3.33) vào (3.39) ta được

$$\begin{aligned} p(x|\mathbf{D}) &= \int \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \frac{1}{\sigma_n\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{\mu-\mu_n}{\sigma_n}\right)^2\right] d\mu \\ &= \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2}\frac{(x-\mu_n)^2}{\sigma^2 + \sigma_n^2}\right] A, \end{aligned} \quad (3.40)$$

trong đó hệ số  $A$  độc lập với  $x$  và được định nghĩa bởi

$$A = \int \exp \left[ -\frac{1}{2} \frac{\sigma^2 + \sigma_n^2}{\sigma^2 \sigma_n^2} \left( \mu - \frac{\sigma_n^2 x + \sigma^2 \mu_n}{\sigma^2 + \sigma_n^2} \right)^2 \right] d\mu. \quad (3.41)$$

Từ (3.40), chúng ta có thể thấy rằng  $p(x|\mathbf{D})$  là một phân phối chuẩn với mean là  $\mu_n$  và phương sai là  $\sigma^2 + \sigma_n^2$ ;  $p(x|\mathbf{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$ . Hơn nữa, chúng ta cũng dễ dàng thấy rằng  $x$  có thể được xem là tổng của hai biến ngẫu nhiên với mật độ phân phối  $p(\mu|\mathbf{D}) \sim N(\mu_n, \sigma_n^2)$  và  $p(w) \sim N(0, \sigma^2)$ , do đó hàm mật độ phân phối của nó có dạng  $p(x|\mathbf{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$ .

Sau khi xác định  $p(x|\mathbf{D})$  cho từng lớp ta được  $p(x|\omega, \mathbf{D})$  cùng với xác suất tiên nghiệm  $P(\omega)$  sẽ có đủ thông tin để xây dựng bộ phân lớp.

### ***b) Xét trường hợp nhiều chiều***

Phần này tổng quát hóa trường hợp một chiều ở trên, nghĩa là ma trận hiệp phương sai  $\Sigma$  được biết nhưng vector bình quân (mean)  $\mu$  không biết. Như phần trước, các phân phối dạng Gauss như sau



$$p(\mathbf{x}|\boldsymbol{\mu}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ và } p(\boldsymbol{\mu}) \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0), \quad (3.42)$$

trong đó  $\boldsymbol{\mu}_0$ ,  $\boldsymbol{\Sigma}_0$ , và  $\boldsymbol{\Sigma}$  giả sử là đã được xác định. Với tập mẫu  $\mathbf{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , tương tự như (3.31) ta có

$$\begin{aligned} p(\boldsymbol{\mu}|\mathbf{D}) &= a \prod_{j=1}^n p(\mathbf{x}_j|\boldsymbol{\mu}) p(\boldsymbol{\mu}) \\ &= a_1 \exp \left[ -\frac{1}{2} \left( \boldsymbol{\mu}^T (n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1}) \boldsymbol{\mu} \right. \right. \\ &\quad \left. \left. - 2\boldsymbol{\mu}^T \left( \boldsymbol{\Sigma}^{-1} \sum_{j=1}^n \mathbf{x}_j + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right) \right) \right], \end{aligned} \quad (3.43)$$

trong đó các thành phần độc lập  $\boldsymbol{\mu}$  được tích hợp trong hệ số  $a_1$ . Rõ ràng (3.43) có dạng hàm mũ của biểu thức bậc hai

$$p(\boldsymbol{\mu}|\mathbf{D}) = a_2 \exp \left[ -\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n) \right]. \quad (3.44)$$

Như vậy  $p(\boldsymbol{\mu}|\mathbf{D})$  có dạng phân phối Gauss dạng  $p(\boldsymbol{\mu}|\mathbf{D}) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ . Đồng nhất (3.43) và (3.44) ta được

$$\boldsymbol{\Sigma}_n^{-1} = n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1} \quad (3.45)$$

và

$$\Sigma_n^{-1} \mu_n = n \Sigma^{-1} \bar{\mathbf{x}}_n + \Sigma_0^{-1} \mu_0, \quad (3.46)$$

trong đó  $\bar{\mathbf{x}}_n$  là bình quân mẫu (sample mean)

$$\bar{\mathbf{x}}_n = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j. \quad (3.47)$$

Áp dụng tính chất

$$(\mathbf{X}^{-1} + \mathbf{Y}^{-1})^{-1} = \mathbf{X}(\mathbf{X} + \mathbf{Y})^{-1}\mathbf{Y} = \mathbf{Y}(\mathbf{X} + \mathbf{Y})^{-1}\mathbf{X} \quad (3.48)$$

vào (3.45) ta được

$$\Sigma_n = \frac{1}{n} \Sigma \left( \frac{1}{n} \Sigma + \Sigma_0 \right)^{-1} \Sigma_0 \quad (3.49)$$

và

$$\mu_n = \Sigma_0 \left( \frac{1}{n} \Sigma + \Sigma_0 \right)^{-1} \bar{\mathbf{x}}_n + \frac{1}{n} \Sigma \left( \frac{1}{n} \Sigma + \Sigma_0 \right)^{-1} \mu_0. \quad (3.50)$$

Nhận thấy rằng  $\mathbf{x}$  được xem như là tổng của hai biến ngẫu nhiên với mật độ phân phối  $p(\boldsymbol{\mu}|\mathbf{D}) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$  và  $p(\mathbf{w}) \sim N(\mathbf{0}, \boldsymbol{\Sigma})$ . Do đó hàm mật độ phân phối  $p(\mathbf{x}|\mathbf{D})$  sẽ là

$$p(\mathbf{x}|\mathbf{D}) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n + \boldsymbol{\Sigma}). \quad (3.51)$$

Ngoài ra, chúng ta cũng có thể xác định  $p(\mathbf{x}|\mathbf{D})$  dựa vào  $p(\mathbf{x}|\mathbf{D}) = \int p(\mathbf{x}|\boldsymbol{\mu}) p(\boldsymbol{\mu}|\mathbf{D}) d\boldsymbol{\mu}$ .

### 3.6 Ước lượng tham số Bayes: Tổng quát

Phần trên giới thiệu phương pháp ước lượng Bayes trong trường hợp phân phối Gauss. Tiếp cận này có thể được tổng quát cho các trường hợp phân phối khác đã được biết dạng của nó. Quá trình được tóm lược như sau:

- Dạng của hàm mật độ xác suất có điều kiện (hay likelihood)  $p(\mathbf{x}|\boldsymbol{\theta})$  được biết, nhưng vector tham số của nó,  $\boldsymbol{\theta}$ , chưa biết và ta cần phải xác định.
- Thông tin về tham số  $\boldsymbol{\theta}$  chỉ được biết qua hàm mật độ xác suất tiên nghiệm  $p(\boldsymbol{\theta})$ .
- Thông tin còn lại liên quan đến  $\boldsymbol{\theta}$  là tập dữ liệu  $\mathbf{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  bao gồm các mẫu được rút trích độc lập từ hàm mật độ phân phối chưa biết  $p(\mathbf{x})$ .

Bài toán trở thành tính mật độ xác suất hậu nghiệm  $p(\theta|D)$ , từ đó chúng ta tính  $p(x|D)$  theo công thức (3.28). Sử dụng công thức Bayes ta có

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}, \quad (3.52)$$

trong đó

$$p(D) = \int p(D|\theta)p(\theta)d\theta \quad (3.53)$$

đóng vai trò là hệ số chuẩn hóa, nó là tích phân đối với tham số  $\theta$ . Trong trường hợp phân phối được giả thiết là phân phối Gauss như trên hoặc dữ liệu với số chiều nhỏ thì việc xác định có thể được thực hiện. Tuy nhiên khi phân phối là bất kỳ và/hoặc dữ liệu với số chiều lớn, việc tính tích phân này trở nên phức tạp hơn. Một số tiếp cận để giải quyết bài toán này bao gồm:

- Sử dụng tích phân số.
- Xấp xỉ các hàm được sử dụng để tính hậu nghiệm bởi các hàm đơn giản hơn (variational Bayes).
- Sử dụng các phương pháp Monte Carlo, cụ thể là Markov Chain Monte Carlo (MCMC).