

# TÌM CHÂN LÝ TỪ DỮ LIỆU THÔ: SUY DIỄN CHỦ ĐỀ VÀ HÌNH HỌC <sup>1</sup>

Ngày nay chúng ta quen dần với sự hiện diện của các thuật toán học máy thông minh. Chúng tự động rà quét qua các nguồn thông tin khổng lồ để truy tìm những nội dung ẩn sau rừng dữ liệu thô sơ và hỗn độn. Một số thuật toán lấy ở đâu vào hàng vạn, triệu trang văn bản tin tức trên mạng, ở đâu ra chúng cho ta biết những bài viết về giáo dục, thể thao, hay nghệ thuật. Hoặc những chủ đề tinh hơn, về thi cử, về doping trong thể thao, về một phong cách mới nổi của nghệ thuật đương đại. Cũng những thuật toán này, khi được áp vào bộ mã di truyền của các cá thể trong một và nhiều quần thể loài người, đã góp phần giúp các nhà khoa học kiến tạo những giả thuyết lý thú về nguồn gốc và sự pha trộn tổ tiên của loài người trong lịch sử xa xôi.

Những thuật toán kể trên dựa vào một nền tảng chung là một lớp mô hình thống kê, gọi là mô hình chủ đề (*topic models*). Đây hiện là một công cụ phổ biến và hữu hiệu trong việc tự động tìm tòi nội dung các văn bản, không chỉ văn bản mà cả nhiều dạng dữ liệu khác, lấy từ kho những bức ảnh chụp, thư viện điện tử, bộ mã di truyền, mạng xã hội. Mô hình thống kê là một công cụ toán học đặc biệt để mô tả motif của dữ liệu: từ đó cho ta các phương pháp tìm tòi suy diễn về các *motif*, hay *pattern* ẩn sau các tần suất về thống kê của dữ liệu thô.

Trong bài viết này tôi sẽ giới thiệu và tập trung vào một vài khía cạnh toán học căn bản của lớp mô hình chủ đề và các phương pháp thống kê kéo theo. Đặc biệt, có một sự liên hệ mật thiết giữa mô hình chủ đề và hình học lồi, giúp cho ta hiểu rõ hơn bản chất của thuật toán học máy đề cập ở trên. Nó cũng giúp ta tạo ra các thuật toán suy diễn mới, hiệu quả hơn.

## 1 Mô hình chủ đề

Mô hình *Phân bố ẩn Dirichlet* (tiếng Anh: *Latent Dirichlet Allocation* (LDA)) được David Blei, Andrew Ng và Michael Jordan thuộc Đại học California, Berkeley giới thiệu vào năm 2001 (Blei et al, 2003). Nó đóng một vai trò đặc biệt quan trọng trong ngành học

máy và trí tuệ nhân tạo. Một mô hình tương tự, về mặt toán học, được phát kiến độc lập bởi một nhóm tác giả khác thuộc ĐH Oxford, Jonathan Pritchard, Matthew Stephens và Peter Donnelly, cũng có một vị trí quan trọng trong ngành di truyền quần thể (Pritchard et al, 2000). Hai bài báo này gộp lại hiện có trên 30,000 trích dẫn trên Google Scholar, cho thấy sự ảnh hưởng đáng kinh ngạc của

<sup>1</sup>Tác giả: Nguyễn Xuân Long, Đại học Michigan, USA. Bài báo đã được đăng trên tạp chí Pi của hội Toán học Việt Nam, số tháng 5, 2017.

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services." Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

Hình 1: Hình bên trái là ví dụ của một trong nhiều văn bản đầu vào của thuật toán. Hình bên phải là đầu ra của thuật toán, cho biết những chủ đề chính của văn bản ở bên trái. Mỗi chủ đề được minh hoạ theo từng cột bằng các từ vựng có tần suất xuất hiện cao (nguồn: Blei et al (2003)).

chúng trong những địa hạt hoàn toàn cách biệt, từ khoa học đến công nghệ. Chỉ với kiến thức toán sơ cấp, chúng ta có thể mô tả kỹ thuật tiên tiến này một cách khá dễ dàng.

Trước khi mô tả LDA, cần có đôi lời về sự cần thiết của mô hình xác suất trong suy diễn thống kê. Lấy  $x$  làm ký hiệu cho đầu vào của một thuật toán, còn  $\theta$  là đầu ra. Một cách cụ thể  $x$  là biểu diễn của tập các văn bản, còn  $\theta$  là biểu diễn cho các chủ đề chúng ta muốn chất lọc từ nội dung của dữ liệu  $x$ .  $x$  lấy giá trị ở trong một không gian mà chúng ta sẽ định nghĩa một cách hình thức. Tương tự như vậy với  $\theta$ . Để có thể suy diễn được chủ đề  $\theta$  từ dữ liệu thô  $x$ , chúng ta cần phải thiết lập được một sự liên hệ giữa  $x$  và  $\theta$ . Sự liên hệ ấy sẽ được mô tả bằng một phân bố xác suất, ký hiệu bởi  $p(x|\theta)$ , mà trong đó  $\theta$  đóng vai trò là một tham số của phân phối cho biến ngẫu nhiên  $x$ .

Trong thảo luận của chúng ta về sự thiết lập của mô hình, dữ liệu được xem là hiện sinh (realization) của một biến ngẫu nhiên  $x$ , còn chủ đề  $\theta$  là một tham số cho phân phối của dữ liệu. Việc sử dụng ngôn ngữ xác suất để xác lập mô hình toán học cho dữ liệu được xem là hết sức tự nhiên, vì dữ liệu thường có yếu tố *bất định*: mỗi lần thu thập dữ liệu mới ta thường nhận một giá trị khác nhau, cho dù quy luật sinh ra dữ liệu có thể được xác định bởi cùng một phân bố "chân lý". Dầu ta không bao giờ biết chắc chắn được phân bố chân lý, bằng những hiểu biết từ dữ liệu thực tế chúng ta có thể tìm một lớp phân bố xấp

xỉ khả dĩ cho phân bố chân lý. Lớp phân bố ấy sẽ được tham số hoá bởi  $\theta$ . Bằng dữ liệu thực, chúng ta sẽ tìm giá trị tốt nhất có thể cho  $\theta$ .

**Mô hình sinh  $p(x|\theta)$**  Bây giờ chúng ta đã sẵn sàng với việc thiết lập mô hình chủ đề một cách chặt chẽ, còn gọi là *mô hình sinh* cho dữ liệu. Đó chẳng qua là phân bố xác suất  $p(x|\theta)$ . Giả sử  $V$  là bộ từ điển gồm có  $d$  phần tử (từ vựng) được đánh số bằng  $V = \{1, \dots, d\}$ . Mỗi văn bản được xem là một hiện sinh (realization) của biến ngẫu nhiên  $x$ . Không mất tính tổng quát, giả dụ văn bản dài  $n$  từ, ta viết  $x = (x_1, \dots, x_n) \in V^n$ .

Giả sử rằng có  $k$  chủ đề ẩn khác nhau có thể dùng để mô tả nội dung tập các văn bản ta có trong tay.  $k$  chủ đề được biểu diễn bằng  $\theta = (\theta_1, \dots, \theta_k)$ , tương ứng với  $k$  phần tử nằm trong đơn hình xác suất

$$\Delta^{d-1} := \{u \in [0, 1]^d | u_1 + \dots + u_d = 1\}.$$

Mỗi chủ đề  $\theta_i \in \Delta^{d-1}$  tương ứng với một tần suất nhất định của các từ vựng trong bộ từ điển. Ví dụ: một chủ đề về "giáo dục" sẽ có nhiều từ với tần suất xuất hiện cao, như "trường", "lớp", "học sinh", "học phí", v.v. Một chủ đề về "nghệ thuật" sẽ có những từ khác với tần suất lớn, như "nghệ sĩ", "mới", "biểu diễn", "tình yêu", v.v. (xem minh hoạ ở Hình 1). Thật khó tưởng tượng một văn bản chỉ đề cập đến một chủ đề duy nhất. Thực tế hơn, mỗi văn bản bao gồm sự pha trộn

của nhiều chủ đề cùng một lúc, dấu mức độ pha trộn có thể nhiều ít khác nhau. (Còn ở ứng dụng trong sinh học phân tử, mỗi cá thể trong một cộng đồng có thể là kết quả của sự pha trộn các nguồn gốc khác nhau về mã di truyền, từ châu Phi, Á hay Âu).

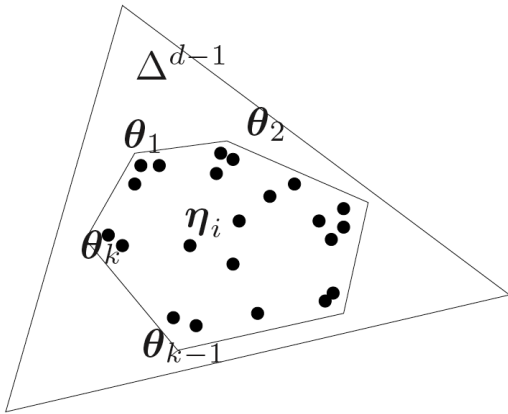
Thực quan trên gợi cho ta một khái niệm hình học rất tự nhiên: điểm ngẫu nhiên nằm trong *tập bao lồi* của các chủ đề:  $G := \text{conv}\{\theta_1, \dots, \theta_k\} \subset \Delta^{d-1}$ . Điểm ngẫu nhiên này, với ký hiệu  $\eta$ , được xác lập bởi công thức

$$\eta = \beta_1 \theta_1 + \dots + \beta_k \theta_k,$$

trong đó  $\beta = (\beta_1, \dots, \beta_k)$  là một biến ngẫu nhiên lấy giá trị trong đơn hình xác suất  $\Delta^{k-1}$ . Chúng ta sẽ giả dụ rằng  $\beta$  tuân theo phân phối Dirichlet trong  $\Delta^{k-1}$ , phân bố này có hàm mật độ như sau, cho mỗi  $\beta \in \Delta^{k-1}$ ,

$$p(\beta|\gamma) = \frac{\Gamma(\sum_{j=1}^k \gamma_j)}{\prod_{j=1}^k \Gamma(\gamma_j)} \prod_{j=1}^k \beta_j^{\gamma_j-1}.$$

Trong công thức trên đây,  $\Gamma$  biểu thị hàm số Gamma, còn  $\gamma_1, \dots, \gamma_k$  là những tham số của mật độ Dirichlet. Định nghĩa trên đây cho (gián tiếp) một phân phối đối với biến ngẫu nhiên  $\eta$  lấy giá trị ở trong đa diện lồi  $G$ . Ta có thể viết ra một cách tường minh phân phối này, ký hiệu bởi  $dP(\eta|\theta, \gamma)$ , bằng công thức chuyển biến (điều này không cần thiết cho bài viết).



Để ý rằng tham số  $\gamma_1, \dots, \gamma_k$  điều khiển mật độ của biến ngẫu nhiên  $\eta \in G$ . Ví dụ, nếu

như  $\gamma_1 = \dots = \gamma_k = 1$ ,  $p(\cdot|\gamma)$  trở thành hằng số. Do đó khi  $k < d$ , nếu  $\theta_1, \dots, \theta_k$  là  $k$  đỉnh của đa diện lồi  $G$ , biến ngẫu nhiên  $\eta$  được định nghĩa như ở trên sẽ tuân theo phân bố *đồng nhất* trong lòng đa diện  $G$ . Nếu  $\gamma_1 = \dots = \gamma_k \ll 1$ , phân phối của  $\eta$  tập trung phần lớn khối (mass) ở sát các mặt biên của  $G$ . Ngược lại nếu  $\gamma_1 = \dots = \gamma_k \gg 1$ , phần lớn khối tập trung ở sâu bên trong lòng đa diện  $G$ .

Chúng ta đã gần xong việc định nghĩa cho phân bố của dữ liệu  $x$ . Mỗi văn bản sẽ có tương ứng một tần suất  $\eta$  kể trên. Khi biết  $\eta$ , văn bản  $x$  được coi là một dãy các từ vựng trong  $V$ . Như đã định ở trên, văn bản  $n$  từ được viết thành  $x = (x_1, \dots, x_n) \in V^n$ . Mỗi  $x_i$  là một biến ngẫu nhiên độc lập tuân theo phân phối phân loại: tung một con xúc xắc  $d$  mặt sao cho xác suất của  $x_i$  lấy giá trị mặt  $j$  sẽ là  $\eta_j$ ,  $j = 1, \dots, d$ . Công thức hàm khối xác suất cho  $x$ , khi biết  $\eta$ :

$$p(x|\eta) = \prod_{i=1}^n \prod_{j=1}^d \eta_j^{1(x_i=j)}. \quad (1)$$

Ở trên,  $1(A) := 1$  nếu biểu hiện  $A$  đúng, và 0 nếu  $A$  sai.

Tóm lại, nếu đã biết giá trị các tham số  $\theta_1, \dots, \theta_k$ , cũng như  $\gamma_1, \dots, \gamma_k$  được cho trước, mô hình sinh dữ liệu được định nghĩa bởi hàm mật độ sau đây

$$p(x|\theta) = \int p(x|\eta) dP(\eta|\theta, \gamma). \quad (2)$$

Ở công thức trên đây,  $dP(\eta|\theta, \gamma)$  là ký hiệu cho phân bố của biến ngẫu nhiên  $\eta$ .

## 2 Suy diễn thống kê đối với chủ đề và hình học lồi

Trên đây chúng ta đã định nghĩa mô hình Phân bố ẩn Dirichlet thông qua phân bố xác suất  $p(x|\theta)$ . Giả sử rằng chúng ta có trong tay  $m$  văn bản, chúng sẽ được xem là  $m$  hiện

sinh độc lập của biến ngẫu nhiên  $\mathbf{x}$ , và được ký hiệu bởi  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ . Mỗi văn bản có đúng  $n$  từ. Như vậy, lượng dữ liệu thô là  $mn$  từ trong  $m$  văn bản. Bài toán thống kê là phải tìm ra các chủ đề  $\theta$  được dùng để xác định phân bố sinh ra tập  $m$  văn bản.

**Bài toán 1.** Cho trước số chủ đề  $k$  và các tham số  $\gamma_1, \dots, \gamma_k > 0$ . Với các văn bản  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$  là  $m$  hiện sinh độc lập của phân phối  $p(\mathbf{x}|\theta)$ , hãy tìm các chủ đề  $\theta = (\theta_1, \dots, \theta_k)$ .

Chúng ta không thể tìm ra chính xác các chủ đề  $\theta$ , chỉ có thể ước lượng chúng. Nhưng ta kỳ vọng rằng nếu số lượng dữ liệu  $m, n$  càng lớn thì các ước lượng đối với  $\theta$  càng tiến gần đến giá trị “chân lý” của chúng. Mặt khác, đây cũng không phải là vấn đề đơn giản về mặt tính toán, mà phải cần những hỗ trợ về thuật toán đủ mạnh: đó là các thuật toán “học máy” lấy đầu vào là tập  $m$  văn bản, và đầu ra chính là các ước lượng về “chủ đề” như tôi đưa ra ở đầu bài viết.

Để hiểu bản chất các thuật toán đó, cần quay về một bài toán dễ hơn trong hình học lồi.

**Ước lượng các đỉnh của một đa diện** Có một bài toán kinh điển trong hình học lồi: làm thế nào để ước lượng được các đỉnh của một đa diện lồi  $G$  nếu như chúng ta quan sát được  $m$  hiện sinh độc lập rút ra theo phân bố đồng nhất từ trong lòng của  $G$ . Bài toán sau đây khái quát hơn một chút, ở chỗ sử dụng một phân bố Dirichlet với một vector tham số  $\gamma = (\gamma_1, \dots, \gamma_k)$  đã được cho trước.

**Bài toán 2.** Cho trước  $\gamma \in \mathbb{R}_+^k$ . Giả sử  $\eta^{(1)}, \dots, \eta^{(m)}$  là  $m$  hiện sinh độc lập của một biến ngẫu nhiên  $\eta$  lấy giá trị trong  $G$  như đã định nghĩa ở mục trước. Hãy đưa ra một ước lượng tốt nhất cho đa diện  $G$ .

“Tốt nhất” được đo bằng một metric cụ thể, như Hausdorff: với hai đa diện  $G$  và  $G'$  bất kỳ trong một không gian chung  $\mathbb{R}^d$

$$d_H(G, G') = \min\{\epsilon \geq 0 | G \subset G'_\epsilon, G' \subset G_\epsilon\}$$

trong đó  $G_\epsilon = \{\theta + \nu | \theta \in G, \nu \in \mathbb{R}^d, \|\nu\|_2 \leq \epsilon\}$ .

Có một ước lượng rất tự nhiên và kinh điển: lấy hình bao lồi của các hiện sinh  $\eta^{(1)}, \dots, \eta^{(m)}$ ,

$$\hat{G} := \text{conv}(\eta^{(1)}, \dots, \eta^{(m)}).$$

Bỏ qua chi tiết  $\hat{G}$  không nhất thiết có  $k$  đỉnh, ước lượng này có một đảm bảo cần thiết: có thể chứng minh được rằng, nếu  $m \rightarrow \infty$ ,  $d_H(G, \hat{G}) \rightarrow 0$  trong xác suất. Kết quả này là một dạng luật số lớn trong hình học lồi. Ngoài ra còn có một số kết quả tinh hơn về tốc độ hội tụ, và về phân bố tiệm cận (Reitzner, 2005; Bárány & Vu, 2007).

Tuy khá giống nhau, bài toán suy diễn về chủ đề (Bài toán 1) của chúng ta phức tạp hơn Bài toán 2 ở một số khía cạnh. Thứ nhất, ở Bài toán 1 ta không có trong tay các hiện sinh của vector tần suất  $\eta^{(1)}, \dots, \eta^{(m)}$ , thay vào đó là các hiện sinh  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$  ( $\mathbf{x}$  liên hệ với  $\eta$  qua phương trình (1)). Thứ hai, vector tham số  $\gamma$  có thể bao gồm các giá trị khác với 1, và khác nhau (để phản ánh xác thực hơn dữ liệu thực tế). Nói cách khác, phân phối trong lòng đa diện lồi  $G$  thường không bao giờ là đồng nhất, tuy đó là một giả thuyết khá phổ biến trong hình học lồi cổ điển.

**Ước lượng bằng cực đại của hàm khả dĩ** Trong thống kê toán có một phương pháp ước lượng rất mạnh và hữu dụng, dựa vào phép lấy cực đại của hàm khả dĩ của Fisher. Hàm khả dĩ (likelihood function) là một hàm số đối với tham số  $\theta$  của (mật độ) phân bố  $p(\mathbf{x}|\theta)$ . Ước lượng cực đại của hàm khả dĩ, viết tắt bằng MLE, khi đã được cho các dữ liệu  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$  được định nghĩa như sau:

$$\hat{\theta} := \arg \max_{\theta} \sum_{i=1}^m \log p(\mathbf{x}^{(i)}|\theta) \quad (3)$$

trong đó hàm khả dĩ  $p(\mathbf{x}^{(i)}|\cdot)$  được cho bởi công thức (2).

Định lý sau cho một bảo đảm quan trọng đối với ước lượng MLE cho lớp mô hình chủ đề:

**Định lý 2.1.** Giả sử  $G$  có  $k$  đỉnh  $\theta_1, \dots, \theta_k$ , và  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$  là ước lượng MLE định nghĩa bởi (3). Viết  $\hat{G} := \text{conv}(\hat{\theta}_1, \dots, \hat{\theta}_k)$ . Nếu  $m$  và  $n$  cùng tiến tới vô cùng, trong đó  $\log \log m \leq \log n = o(m)$ , ta có

$$d_H(\hat{G}, G) = O_p \left( \frac{\log m}{m} + \frac{\log n}{n} + \frac{\log n}{m} \right)^{1/(2(q+\alpha))}$$

Trên đây,  $q = \min\{(k-1), (d-1)\}$ ,  $\alpha$  là một hằng số dương phụ thuộc vào  $\gamma$ ,  $O_p$  có nghĩa là “hằng số cận trên trong xác suất”, được xác định đối với phân phối  $p(\mathbf{x}|\theta)$ .

Với độc giả chưa có đủ kiến thức về lý thuyết xác suất để hiểu kết quả trên một cách cặn kẽ, chỉ cần biết rằng ước lượng MLE (3) càng tiến gần tới chủ đề chân lý khi lượng dữ liệu càng nhiều. Định lý cho ta biết chặn trên của sai số của ước lượng MLE, khi  $m$  và  $n$  tăng dần. Hãy tham khảo Nguyen (2015) để biết thêm chi tiết cho cả chặn trên và chặn dưới.

**Bayes hay không Bayes** Chúng ta tạm dùng một chốc lát, để quay lại với thảo luận về sự cần thiết của mô hình xác suất  $p(\mathbf{x}|\theta)$ , một chất keo toán học kết nối dữ liệu quan sát được  $\mathbf{x}$ , với khái niệm “chân lý” ẩn sau đó,  $\theta$ . Các nhà thống kê đều đồng ý với nhau rằng dữ liệu phải được xem là hiện sinh của một biến ngẫu nhiên, do sự bất định của quan sát. Nhưng có một sự tranh cãi quyết liệt về ý nghĩa của khái niệm  $\theta$ : nó có bất định hay không? Có hai trường phái đối với câu hỏi này: Các nhà *thống kê tần suất* cho rằng  $\theta$  hoàn toàn có thể xác định được. Phương pháp ước lượng cực đại của hàm khả dĩ (MLE) chính là một sản phẩm của các nhà tần suất. Ngược lại, trường phái *Bayes* cho rằng tri thức về  $\theta$  cũng bất định, và do đó về mặt toán học, nó cũng phải được biểu diễn bởi một biến ngẫu nhiên. Vì chúng ta không thể quan sát được khái niệm trừu tượng  $\theta$ , chúng ta không thể xác quyết đúng sai trong tranh cãi này.

Một ưu thế của phương pháp suy diễn Bayes là nó không chỉ cung cấp cho ta một điểm ước lượng cho  $\theta$ , nó còn cho ta biết mức độ

(không) chắc chắn của ước lượng đó, dựa trên những dữ liệu ta thu thập được. Để áp dụng phương pháp Bayes, tham số  $\theta$ , bởi được xem là biến ngẫu nhiên, phải được gán cho một phân bố trước khi có dữ liệu. Phân phối này được gọi là *phân bố tiên nghiệm*, ký hiệu bởi  $\Pi(\theta)$ , nó là chất lọc của tri thức mà nhà thống kê có được trước khi anh ta quan sát. Sau khi đã thu thập dữ liệu,  $m$  văn bản  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ , kiến thức của nhà thống kê về  $\theta$  được tổng kết thành *phân bố hậu nghiệm*. Phân bố hậu nghiệm được xác định bằng công thức Bayes nổi tiếng:

$$d\Pi(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)d\Pi(\theta)}{\int p(\mathbf{x}|\theta)d\Pi(\theta)}. \quad (4)$$

Lưu ý, cả phân bố tiên nghiệm (a priori) và hậu nghiệm (a posteriori) thường được ký hiệu bởi chữ  $\Pi$  viết hoa, để phân biệt với hằng số  $\pi$  thân thuộc với độc giả tạp chí này. Có thể chứng minh được rằng khi số lượng dữ liệu càng nhiều thì phân bố hậu nghiệm tập trung càng gần về giá trị “chân lý” của  $\theta$  (Nguyen, 2015): niềm tin của các nhà thống kê Bayes cũng tiệm cận dần tới chân lý!

### 3 Thuật toán học chủ đề

Thống kê hiện đại và học máy tập trung đáng kể vào mục tiêu tính toán hiệu quả những đại lượng trọng yếu như ước lượng MLE (3) hay phân bố hậu nghiệm (4). Đây là một thách thức vì những phép tính này nằm trong lớp NP-hard hoặc khó hơn. Trên thực tế, để có độ tin cậy cao các thuật toán phải quét qua một số lớn các văn bản có độ dài khác nhau, con số đó có thể đến hàng vạn, hàng triệu, thậm chí hàng tỷ, như Định lý 2.1 gợi ý. Khi đó, thời gian chạy thuật toán có thể mất hàng ngày hoặc hàng tuần để có được kết quả tương đối chính xác của bài toán tìm cực đại cho hàm khả dĩ.

Vì thế, các nhà thống kê và học máy thực hành chỉ cần tìm các giải pháp xấp xỉ, như phép suy diễn biến phân (variational inference) (Blei et al, 2003). Đối với phương

pháp suy diễn Bayes, phân bố hậu nghiệm cũng được xấp xỉ bằng cách thiết kế một thuật toán mô phỏng chuỗi Markov với thuộc tính ergodic, mà theo đó phân bố dừng của nó chính là phân bố hậu nghiệm cần phải tính. Có một thực trạng khó chịu: phép biến phân xấp xỉ cho ra kết quả rất nhanh, nhưng lại không đúng. Ngược lại, thuật toán mô phỏng chuỗi Markov đưa ra được ước lượng khá đúng, nhưng lại không nhanh.

Sự liên hệ giữa mô hình chủ đề và hình học lồi được mô tả trên đây đã được tận dụng để đưa ra một thuật toán ước lượng chủ đề vừa chính xác, vừa nhanh, trong một số điều kiện nhất định. Vì giới hạn bài vở, tôi chỉ mô tả một cách nôm na ý tưởng này như sau: Giả sử chúng ta biết rằng phần lớn các văn bản tập trung ở gần các đỉnh của đa diện  $G$  (nếu như tham số  $\gamma_1, \dots, \gamma_k \ll 1$ ). Ta có thể áp dụng một thuật toán phân cụm đơn giản để tìm ra trọng tâm của  $k$  cụm các văn bản này. Các trọng tâm này cho ta một ước lượng thô khả dĩ về vị trí các đỉnh của  $G$ . Để có ước lượng tốt hơn ta cần một số phép “điều chỉnh hình học”. Độc giả tò mò có thể tham khảo thêm Yurochkin và Nguyen (2016).

## 4 Kết luận

Khoa học dữ liệu phát triển mạnh mẽ như hôm nay là nhờ đóng góp về vật chất của công nghệ thông tin, bên cạnh vai trò mang tính nền tảng của thống kê, xác suất nói riêng và toán học nói chung. Công nghệ thông tin mang cho ta nguồn dữ liệu dồi dào, làm nảy sinh một nhu cầu cần suy diễn với những dữ liệu ấy. Thống kê toán cho ta một nền móng để suy diễn một cách hợp lý. Ngược lại sự phát triển của khoa học suy diễn từ dữ liệu, dữ liệu lớn, cũng đưa toán học tới gần với các ứng dụng của xã hội hiện đại hơn bao giờ hết. Vấn đề suy diễn với các cấu trúc phức tạp cũng có tác dụng thúc đẩy sự phát triển nội tại của toán học.

Bài viết này minh họa một vài khía cạnh của sự tương tác ấy: chỉ với kiến thức toán sơ cấp chúng ta có thể mô tả được một lớp kỹ thuật tiên tiến của ngành học máy và di truyền quần thể trong hai thập niên đầu thế kỷ 21. Tất nhiên, để sáng tạo ra những kỹ thuật ấy, các tác giả của chúng không chỉ cần có toán học. Họ phải có sự tò mò đối với thực tiễn, có khả năng diễn đạt thành thực bằng mô hình toán học, có sự liêu lĩnh để áp mô hình vào dữ liệu thực. Để hiểu thấu đáo và cải thiện các phương pháp suy diễn ấy, để đi đến những bước đột phá trong tương lai, chắc chắn sẽ cần nhiều công cụ toán học mạnh hơn nữa.

**Lời cảm ơn** Tác giả chân thành cảm ơn GS Ngô Bảo Châu đã khích lệ và góp ý cho bài viết.

## 5 Tài liệu tham khảo

- I. Bárány and V. Vu. Central limit theorems for Gaussian polytopes. *Annals of Probability*, 35:1593–1621, 2007.
- D. Blei, A. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- X. Nguyen. Posterior contraction of the population polytope in finite admixture models. *Bernoulli*, 21:618–646, 2015.
- J. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multilocus genotype data. *Genetics*, 155:945–959, 2000.
- M. Reitzner. Central limit theorems for random polytopes. *Probability Theory and Related Fields*, 133:483–507, 2005.
- M. Yurochkin and X. Nguyen. Geometric Dirichlet means algorithm for topic inference. *Advances in Neural Information Processing Systems (NIPS)* 29:2505–2513, 2016.