

Probability & Statistics – Xác Suất & Thống Kê

Nguyễn Quân Bá Hồng*

Ngày 12 tháng 10 năm 2024

Tóm tắt nội dung

This text is a part of the series *Some Topics in Advanced STEM & Beyond*:

URL: https://nqbh.github.io/advanced_STEM/.

Latest version:

- *Probability & Statistics – Xác Suất & Thống Kê*.

PDF: URL: https://github.com/NQBH/advanced_STEM_beyond/blob/main/probability_statistics/NQBH_probability_statistics.pdf.

TEX: URL: https://github.com/NQBH/advanced_STEM_beyond/blob/main/probability_statistics/NQBH_probability_statistics.tex.

Mục lục

1 Basic	1
2 Data Science (DS)	1
3 Deep Learning (DL)	1
4 Machine Learning (ML)	2
5 Artificial Intelligence (AI)	2
6 Miscellaneous	3
Tài liệu	3

1 Basic

Relationship among AL, ML, & DL. $DL \subset ML \subset AI$.

2 Data Science (DS)

3 Deep Learning (DL)

Resources – Tài nguyên.

1. [LBH15]. YANN LECUN, YOSHUA BENGIO, GEOFFREY HINTON. *Deep Learning*.

Những năm gần đây, sự phát triển của các hệ thống tính toán cùng lượng dữ liệu khổng lồ được thu thập bởi các hãng công nghệ lớn đã giúp machine learning tiến thêm 1 bước dài. 1 lĩnh vực mới được ra đời được gọi là *học sâu* (deep learning, DL). Deep learning đã giúp máy tính thực thi những việc vào 10 năm trước tưởng chừng là không thể: phân loại cả ngàn vật thể khác nhau trong các bức ảnh, tự tạo chú thích cho ảnh, bắt chước giọng nói & chữ viết, giao tiếp với con người, chuyển đổi ngôn ngữ, hay thậm chí cả sáng tác văn thơ & âm nhạc.” – [Tiệ25, p. 15]

Φ: Start with simple things – Luôn bắt đầu từ những điều đơn giản. Khi bắt tay vào giải quyết 1 bài toán ML hay bất cứ bài toán nào, nên bắt đầu từ các thuật toán đơn giản. Không phải chỉ có các thuật toán phức tạp mới có thể giải quyết được vấn đề. Các thuật toán phức tạp thường có yêu cầu cao về khả năng tính toán & đôi khi nhạy cảm với cách chọn tham số. Ngược lại, các thuật toán đơn giản giúp ta nhanh chóng có 1 bộ khung cho mỗi bài toán. Kết quả của các thuật toán đơn giản cũng mang lại cái nhìn sơ bộ về sự phức tạp của mỗi bài toán. Việc cải thiện kết quả sẽ được thực hiện dần ở các bước sau. ” – [Tiệ25, p. 17]

*A Scientist & Creative Artist Wannabe. E-mail: nguyenquanbahong@gmail.com. Bến Tre City, Việt Nam.

Approach. Để giải quyết mỗi bài toán ML, cần chọn 1 mô hình phù hợp. Mô hình này được mô tả bởi bộ các tham số ta cần đi tìm. Thông thường, lượng tham số có thể lên tới hàng triệu & được tìm bằng cách giải 1 bài toán tối ưu. Khi viết về các thuật toán ML, VKTiếp sẽ bắt đầu từ các ý tưởng trực quan. Các ý tưởng này được mô hình hóa dưới dạng 1 bài toán tối ưu. Các suy luận toán học & ví dụ mẫu trên Python sẽ giúp hiểu rõ hơn về nguồn gốc, ý nghĩa, & cách sử dụng mỗi thuật toán. Xen kẽ giữa các thuật toán ML, trình bày các kỹ thuật tối ưu cơ bản, với hy vọng giúp hiểu rõ hơn bản chất của vấn đề.

Audiences. Cuốn sách được thực hiện hướng tới nhiều nhóm độc giả khác nhau. Nếu không thực sự muốn đi sâu vào phần toán, vẫn có thể tham khảo mã nguồn & cách sử dụng các thư viện. Nhưng để sử dụng các thư viện 1 cách hiệu quả, cũng cần hiểu nguồn gốc của mô hình & ý nghĩa của các tham số. Còn nếu thực sự muốn tìm hiểu nguồn gốc, ý nghĩa của các thuật toán, có thể học được nhiều điều từ cách xây dựng & tối ưu các mô hình.

Python. Python là 1 ngôn ngữ lập trình miễn phí, có thể được cài đặt dễ dàng trên các nền tảng hệ điều hành khác nhau. Có rất nhiều thư viện hỗ trợ ML cũng như DL trên Python. Có 2 thư viện Python chính thường được sử dụng là `numpy`, `scikit-learn`.

- `numpy` www.numpy.org là 1 thư viện phổ biến giúp xử lý các phép toán liên quan đến các mảng nhiều chiều, hỗ trợ các hàm gần gũi với đại số tuyến tính. Cách xử lý các mảng nhiều chiều.
- `scikit-learn/sklearn` scikit-learn.org: 1 thư viện chứa đầy đủ các thuật toán ML cơ bản & rất dễ sử dụng. Tài liệu của `scikit-learn` cũng là 1 nguồn tham khảo chất lượng cho MLer. `Scikit-learn` được dùng để kiểm chứng các suy luận toán học & các mô hình được xây dựng thông qua `numpy`.

Inevitability of mathematics in ML. Có rất nhiều thư viện giúp tạo ra các sản phẩm ML/DL mà không yêu cầu nhiều kiến thức toán. Hướng tới việc giúp hiểu bản chất toán học đằng sau mỗi mô hình trước khi áp dụng các thư viện sẵn có. Việc sử dụng thư viện + yêu cầu kiến thức nhất định về việc lựa chọn mô hình & điều chỉnh các tham số.

4 Machine Learning (ML)

Resources – Tài nguyên.

1. Machine Learning Mastery: Making Developers Awesome at Machine Learning: <https://machinelearningmastery.com>.
 - [Machine Learning Mastery/8 Inspirational Applications of Deep Learning](#).
2. Machine Learning cơ bản: <https://machinelearningcoban.com/>.
3. [Tiếp25]. VŨ HỮU TIẾP. *Machine Learning Cơ Bản*.
Mã nguồn cuốn ebook “Machine Learning Cơ Bản”: <https://github.com/tiepvupsu/ebookMLCB>.

Definition 1. “Machine learning (ML) is a field of study in AI concerned with the development & study of *statistical algorithms* that can learn from *data* & generalize to unseen data, & thus perform *tasks* without explicit *instructions*. Quick progress in the fields of *deep learning*, beginning in 2010s, allowed neural networks to surpass many previous approaches in performance.” – [Wikipedia/machine learning](#)

Định nghĩa 1. “Học máy (*machine learning, ML*) là 1 tập con của trí tuệ nhân tạo. Machine learning là 1 lĩnh vực nhỏ trong Khoa học Máy tính, có khả năng tự học hỏi dựa trên dữ liệu được đưa vào mà không cần phải được lập trình cụ thể: “Machine Learning is the subfield of computer science, that “gives computers the ability to learn without being explicitly programmed” – [Wikipedia](#).” – [Tiếp25, p. 15]

“ML finds application in many fields, including *natural language processing*, *computer vision*, *speech recognition*, *email filtering*, *agriculture*, & *medicine*. The application of ML to business problems is known as *predictive analysis*.”

Statistics & mathematical optimization/mathematical programming methods comprise the foundations of machine learning. *Data mining* is related field of study, focusing on *exploratory data analysis* (EDA) via *unsupervised learning*.

From a theoretical viewpoint, *probably approximately correct (PAC) learning* provides a framework for describing machine learning.” – [Wikipedia/machine learning](#)

Relationships of ML to AI. As a scientific endeavor, machine learning grew out of the quest for AI. In the early days of AI as an *academic discipline*, some researchers were interested in having machines learn from data. They attempted to approach the problem with various symbolic methods, as well as what were then termed “*neural networks*”; these were mostly *perceptrons* & other models e.g. *ADALINE* that were later found to be reinventions of the *generalized linear models* of statistics. *Probabilistic reasoning* was also employed, especially in *automated medical diagnosis*. However, an increasing emphasis on the *logical, knowledge-based approach* caused a rift between AI & machine learning. Probabilistic systems were plagued by theoretical & practical problems of data acquisition & representation.

5 Artificial Intelligence (AI)

Resources – Tài nguyên.

1. [BV14]. LÊ HOÀI BẮC, TÔ HOÀI VIỆT. *Cơ Sở Trí Tuệ Nhân Tạo*.

2. [Aou14]. JOSEPH E. AOUN. *Robot-Proof: Higher Education in the Age of Artificial Intelligence*.
3. [Aou19]. JOSEPH E. AOUN. *Robot-Proof: Higher Education in the Age of Artificial Intelligence – Chạy Dua Với Robot: Học Tập Thời Trí Tuệ Nhân Tạo*.

6 Miscellaneous

Tài liệu

- [Aou14] Joseph E. Aoun. *Robot-Proof: Higher Education in the Age of Artificial Intelligence*. MIT Publisher, 2014, p. 187.
- [Aou19] Joseph E. Aoun. *Robot-Proof: Higher Education in the Age of Artificial Intelligence – Chạy Dua Với Robot: Học Tập Thời Trí Tuệ Nhân Tạo*. Trịnh Huy Nam dịch. Nhà Xuất Bản Thế Giới, 2019, p. 241.
- [BV14] Lê Hoài Bắc and Tô Hoài Việt. *Cơ Sở Trí Tuệ Nhân Tạo*. Nhà Xuất Bản Khoa Học & Kỹ Thuật, 2014, p. 229.
- [LBH15] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep Learning”. In: *Nature* 521 (2015), pp. 436–444. DOI: [10.1038/nature14539](https://doi.org/10.1038/nature14539). URL: <https://doi.org/10.1038/nature14539>.
- [Tiệ25] Vũ Khắc Tiệp. *Machine Learning Cơ Bản*. 2025, p. 422.