

LINEAR REGRESSION

Lê Bích Phương

Trường Đại học Mở - Địa chất

lebichphuong@humg.edu.vn

Tháng 1 năm 2020

1 Lời mở đầu

2 Nội dung

- Problem Formulation
- Parameter Estimation
 - Ước lượng hợp lý cực đại
 - Overfitting in Linear Regression
 - Maximum A Posteriori Estimation

3 Kết luận

4 Tài liệu tham khảo

Tài liệu tham khảo

Hồi quy là một vấn đề cơ bản trong machine learning, các vấn đề hồi quy xuất hiện trong một loạt các lĩnh vực nghiên cứu và ứng dụng, bao gồm phân tích chuỗi thời gian (ví dụ: nhận dạng hệ thống), tối ưu hóa, và các ứng dụng deeplearning (ví dụ: trò chơi máy tính, nhận dạng hình ảnh, chú thích video tự động). Hồi quy cũng là một thành phần then chốt của các thuật toán phân loại. Việc tìm một hàm hồi quy đòi hỏi phải giải quyết nhiều vấn đề khác nhau, bao gồm:

1. **Lựa chọn loại mô hình và tham số** của hàm hồi quy. Khi cho một tập dữ liệu, thì các lớp hàm nào (ví dụ: đa thức, ...) và các tham số nào (ví dụ: bậc của đa thức) ta nên chọn để mô hình hóa dữ liệu.
2. **Tìm các tham số tốt**. Khi đã chọn được mô hình của hàm hồi quy, ta cần tìm các tham số tốt cho mô hình. Ở đây, ta cần chọn tham số để cực tiểu hóa hàm loss.
3. **Overfitting và lựa chọn mô hình**. Overfitting là vấn đề mà hàm hồi quy phù hợp với dữ liệu huấn luyện *quá tốt* nhưng lại không có tính khái quát đối với dữ liệu không nhìn thấy.

4. **Mối quan hệ giữa hàm loss và các tham số priors.** Hàm loss (mục tiêu tối ưu) thường được tạo ra bởi một mô hình xác suất.

Chúng ta sẽ xem xét sự liên quan giữa hàm loss và các giả thuyết cơ bản trước đó gây ra những mất mát này.

5. **Mô hình không chắc chắn.** Trong bất kỳ cài đặt thực tế nào, chúng ta chỉ có quyền truy cập vào một lượng dữ liệu (huấn luyện) hữu hạn, có khả năng lớn để chọn lớp mô hình và các tham số tương ứng. Điều này có nghĩa là lượng dữ liệu huấn luyện hữu hạn thì không bao gồm tất cả các tình huống có thể xảy ra, chúng ta có thể muốn mô tả các tham số không chắc chắn còn lại để đạt được mức độ tin cậy của mô hình dự đoán tại thời điểm thử nghiệm; tập huấn luyện càng nhỏ thì mô hình không chắc chắn càng quan trọng.

Bài toán hồi quy tuyến tính

Ta xét bài toán hồi quy với hàm likelihood như sau:

$$p(y | \mathbf{x}) = \mathcal{N}(y | f(\mathbf{x}), \sigma^2). \quad (1)$$

Trong đó $\mathbf{x} \in \mathbb{R}^D$ là các biến đầu vào và $y \in \mathbb{R}$ là các giá trị hàm mục tiêu. Trong công thức (1), mỗi quan hệ hàm giữa \mathbf{x} và y được cho bởi

$$y = f(\mathbf{x}) + \epsilon, \quad (2)$$

Trong đó $\epsilon = \mathcal{N}(0, \sigma^2)$ là phân phối của biến ngẫu nhiên độc lập, có phân phối Gaussian với giá trị trung bình bằng 0 và phương sai là σ^2 . Mục tiêu của ta là đi tìm một hàm tương tự với hàm f chưa biết.

Ta chọn một hàm được tham số hóa và tìm các tham số θ mà nó “làm tốt” đối với mô hình dữ liệu.

Giả sử rằng phương sai (the noise variance) σ^2 đã biết và tập trung vào việc nghiên cứu các tham số mô hình θ . Trong hồi quy tuyến tính, ta xét trường hợp đặc biệt là các tham số θ xuất hiện tuyến tính trong mô hình.

Bài toán hồi quy tuyến tính

Một ví dụ về hồi quy tuyến tính được đưa cho bởi

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \mathbf{x}^\top \boldsymbol{\theta}, \sigma^2) \quad (3)$$

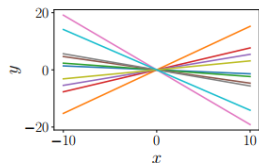
$$\iff y = \mathbf{x}^\top \boldsymbol{\theta} + \epsilon, \quad \epsilon = \mathcal{N}(0, \sigma^2), \quad (4)$$

Trong đó $\boldsymbol{\theta} \in \mathbb{R}^D$ là các tham số cần tìm. Lớp các hàm được mô tả bởi công thức (4) là các đường thẳng đi qua gốc tọa độ. Trong (4), ta đã chọn tham số hóa $f(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\theta}$.

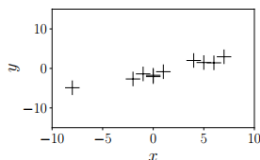
Hàm *likelihood* trong (3) là hàm mật độ xác suất của y được tính tại $\mathbf{x}^\top \boldsymbol{\theta}$.

Example

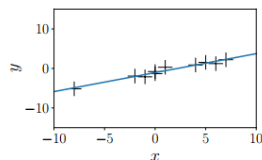
Với $x, \theta \in \mathbb{R}$, mô hình hồi quy tuyến tính trong hình sau mô tả các đường thẳng (các hàm tuyến tính) và tham số θ là hệ số góc của đường thẳng.



(a) Example functions (straight lines) that can be described using the linear model in (9.4).



(b) Training set.

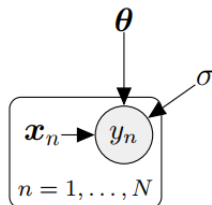


(c) Maximum likelihood estimate.

Hình: Linear regression example. (a) Example functions that fall into this category; (b) training set; (c) maximum likelihood estimate

Ước lượng tham số

Xét cài đặt hồi quy tuyến tính như trong (4) và giả sử cho trước một *tập huấn luyện* (training set) $\mathcal{D} := \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ bao gồm N đầu vào $\mathbf{x}_n \in \mathbb{R}^D$ và các mục tiêu tương ứng $y_n \in \mathbb{R}$, $n = 1, \dots, N$. Mô hình đồ họa tương ứng được thể hiện trong hình sau.



Hình: Probabilistic graphical model for linear regression. Observed random variables are shaded, deterministic/ known values are without circles

Lưu ý rằng y_i và y_j độc lập có điều kiện với các đầu vào tương ứng \mathbf{x}_i , \mathbf{x}_j cho nên

$$p(\mathcal{Y} | \mathbf{X}, \boldsymbol{\theta}) = p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}) \quad (5a)$$

$$= \prod_{n=1}^N p(y_n | \mathbf{x}_n, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{x}_n^\top \boldsymbol{\theta}, \sigma^2), \quad (5b)$$

trong đó ta định nghĩa $\mathbf{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ và $\mathcal{Y} := \{y_1, \dots, y_N\}$ là các tập huấn luyện đầu vào và các mục tiêu tương ứng. The likelihood và các thành phần $p(y_n | \mathbf{x}_n, \boldsymbol{\theta})$ theo phân phối Gaussian.; see (3).

Sau đây, ta sẽ thảo luận về cách tìm ra các tham số tối ưu $\theta^* \in \mathbb{R}^D$ cho mô hình hồi quy tuyến tính (4). Khi các tham số θ^* được tìm ra, ta có thể dự đoán các giá trị hàm theo như công thức (4), tại một giá trị \mathbf{x}_* tùy ý, phân phối của mục tiêu y_* tương ứng là

$$p(y_* | \mathbf{x}_*, \theta^*) = \mathcal{N}(y_* | \mathbf{x}_*^\top \theta^*, \sigma^2). \quad (6)$$

Tiếp theo, ta sẽ xem xét việc đánh giá tham số bằng cách cực đại hóa hàm likelihood.

Ước lượng hợp lý cực đại

Để tìm ra các tham số tối ưu θ_{ML} của bài toán hồi quy tuyến tính, chúng ta cực tiểu hàm đối của log hàm hợp lý

$$-\log p(\mathcal{Y} | \mathcal{X}, \theta) = -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n, \theta) = -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \theta), \quad (7)$$

Trong mô hình hồi quy tuyến tính (4), hàm hợp lý, the likelihood, là Gaussian, do đó ta có

$$\log p(y_n | \mathbf{x}_n, \theta) = -\frac{1}{2\sigma^2} (y_n - \mathbf{x}_n^\top \theta)^2 + \text{const}, \quad (8)$$

trong đó hằng số bao gồm tất cả các thành phần độc lập với θ .

$$\mathcal{L}(\boldsymbol{\theta}) := \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^\top \boldsymbol{\theta})^2 \quad (9a)$$

$$= \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2, \quad (9b)$$

- $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{N \times D}$ như bộ các đầu vào huấn luyện;
- $\mathbf{y} := [y_1, \dots, y_N]^\top \in \mathbb{R}^N$ là vector tất các các mục tiêu huấn luyện

Chú ý rằng hàng thứ n trong ma trận thiết kế, the design matrix, tương ứng với \mathbf{x}_n của đầu vào huấn luyện. Trong (9b), ta đã sử dụng tổng bình phương các lỗi giữa quan sát y_n và giá trị mà mô hình dự đoán tương ứng $\mathbf{x}_n^\top \boldsymbol{\theta}$ bằng bình phương khoảng cách giữa \mathbf{y} và $\mathbf{X}\boldsymbol{\theta}$.

Ta tính gradient của \mathcal{L} theo tham số $\boldsymbol{\theta}$

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \frac{d}{d\boldsymbol{\theta}} \left(\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right) \quad (10a)$$

$$= \frac{1}{2\sigma^2} \frac{d}{d\boldsymbol{\theta}} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \right) \quad (10b)$$

$$= \frac{1}{\sigma^2} (-\mathbf{y}^\top \mathbf{X} + \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}) \in \mathbb{R}^{1 \times D}. \quad (10c)$$

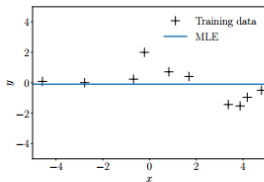
$\boldsymbol{\theta}_{\text{ML}}$ là nghiệm của phương trình $\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \mathbf{0}^\top$

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \mathbf{0}^\top \stackrel{(10c)}{\iff} \boldsymbol{\theta}_{\text{ML}}^\top \mathbf{X}^\top \mathbf{X} = \mathbf{y}^\top \mathbf{X} \quad (11a)$$

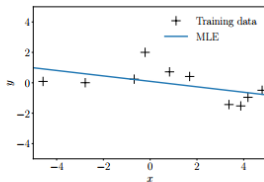
$$\iff \boldsymbol{\theta}_{\text{ML}}^\top = \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (11b)$$

$$\iff \boldsymbol{\theta}_{\text{ML}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (11c)$$

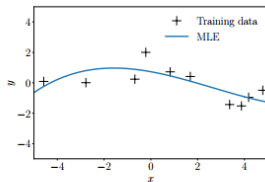
Vấn đề quá mức trong hồi quy tuyến tính



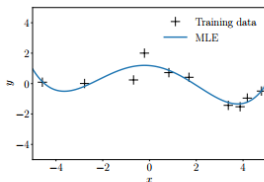
(a) $M = 0$



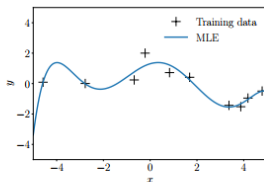
(b) $M = 1$



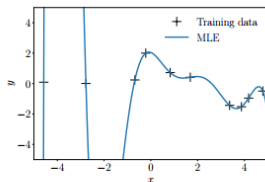
(c) $M = 3$



(d) $M = 4$



(e) $M = 6$



(f) $M = 9$

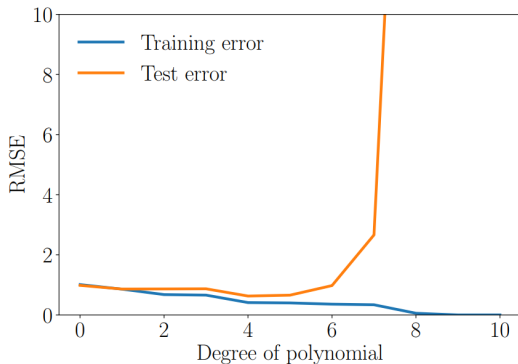
Hình: Maximum likelihood fits for different polynomial degrees M .

Vấn đề quá mức trong hồi quy tuyến tính

Hình trên minh họa một số đa thức ứng với $N = 10$ quan sát.

- Đa thức bậc thấp (hoặc hằng số khi $(M = 0)$ hoặc tuyến tính khi $(M = 1)$) biểu diễn kém,
- Với bậc $M = 3, \dots, 5$ biểu diễn dữ liệu trơn chu hơn,
- $M = N - 1 = 9$, hàm sẽ đi qua mọi điểm của dữ liệu. Tuy nhiên các đa thức bậc cao này dao động mạnh và là sự đại diện kém cho hàm khởi tạo, do vậy ta đối mặt với hiện tượng *overfitting*.

Vấn đề quá mức trong hồi quy tuyến tính



Hình: Training and test error

Bây giờ nhìn vào the test error, ta thấy rằng ban đầu lỗi giảm sau đó tăng vọt lên, và bậc của đa thức bằng 4 là tốt nhất theo nghĩa lỗi của cả Training và Test đều thấp.

Ước lượng hậu nghiệm cực đại

Ước lượng hợp lý cực đại có hiện tượng overfitting. Để giảm thiểu ảnh hưởng của các giá trị tham số quá lớn, chúng ta có thể thay thế tham số θ bởi prior distribution $p(\theta)$ trên các tham số.

Trên một tập dữ liệu \mathcal{X} , \mathcal{Y} có sẵn, thay vì cực đại hóa hàm likelihood ta tìm các tham số mà cực đại hóa the posterior distribution $p(\theta | \mathcal{X}, \mathcal{Y})$.

Việc này được gọi là ước lượng *maximum a posteriori* (MAP).

Ta bắt đầu với dạng log và tính the log-posterior như sau

$$\log p(\theta | \mathcal{X}, \mathcal{Y}) = \log p(\mathcal{Y} | \mathcal{X}, \theta) + \log p(\theta) + \text{const}, \quad (12)$$

ở đó const là số hạng độc lập với θ .

Để tìm ước lượng MPA θ_{MAP} , ta đi cực tiểu hóa hàm đối của log của posterior distribution với biến số θ , nghĩa là ta giải

$$\theta_{\text{MAP}} \in \arg \min_{\theta} \{-\log p(\mathcal{Y} | \mathcal{X}, \theta) - \log p(\theta)\}. \quad (13)$$

Ước lượng hậu nghiệm cực đại

The gradient của hàm đối của log posterior với biến θ là

$$-\frac{d \log p(\theta | \mathcal{X}, \mathcal{Y})}{d\theta} = -\frac{d \log p(\mathcal{Y}, \mathcal{X}, \theta)}{d\theta} - \frac{d \log p(\theta)}{d\theta}, \quad (14)$$

$$p(\theta) = \mathcal{N}(\mathbf{0}, b^2 I)$$

$$-\log p(\theta | \mathcal{X}, \mathcal{Y}) = \frac{1}{2\sigma^2} (\mathbf{y} - \Phi\theta)^\top (\mathbf{y} - \Phi\theta) + \frac{1}{2b^2} \theta^\top \theta + \text{const.} \quad (15)$$

$$-\frac{d \log p(\theta | \mathcal{X}, \mathcal{Y})}{d\theta} = \frac{1}{\sigma^2} (\theta^\top \Phi^\top \Phi - \mathbf{y}^\top \Phi) + \frac{1}{b^2} \theta^\top. \quad (16)$$

$$\frac{1}{\sigma^2}(\boldsymbol{\theta}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} - \mathbf{y}^\top \boldsymbol{\Phi}) + \frac{1}{b^2} \boldsymbol{\theta}^\top = \mathbf{0}^\top \quad (17a)$$

$$\iff \boldsymbol{\theta}^\top \left(\frac{1}{\sigma^2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{1}{b^2} \mathbf{I} \right) - \frac{1}{\sigma^2} \mathbf{y}^\top \boldsymbol{\Phi} = \mathbf{0}^\top \quad (17b)$$

$$\iff \boldsymbol{\theta}^\top \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right) = \mathbf{y}^\top \boldsymbol{\Phi} \quad (17c)$$









$$\iff \boldsymbol{\theta}^\top = \mathbf{y}^\top \boldsymbol{\Phi} \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \quad (17d)$$

vì vậy ước lượng MAP là

$$\boldsymbol{\theta}_{\text{MAP}} = \left(\boldsymbol{\Phi}^\top \boldsymbol{\Phi} + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \boldsymbol{\Phi}^\top \mathbf{y}. \quad (18)$$

Trong bài này, tôi trình bày về vấn đề hồi quy tuyến tính, một trong bốn trụ cột của học máy.

Tài liệu tham khảo

-  C. Cortes, V. Vapnik, in *Machine Learning*, pp. 273–297 (1995).
-  N. Cristianini, J. Shawe-Taylor, *An introduction to support Vector Machines: and other kernel-based learning methods*. Cambridge University Press, New York, NY, USA (2000).
-  J. Shawe-Taylor, N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge University Press, New York, NY, USA (2004).
-  C. J. C. Burges, *Data Mining and Knowledge Discovery* **2**, 121 (1998).
-  C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer (2006).
-  Ian Goodfellow, Yoshua Bengio, and Aaron Courville, *Deep Learning*. 2016
-  Hal Daumé III, *A Course In Machine Learning*, 2018
-  Marc Peter Deisenroth, A.Aldo Faisal, Cheng Soon Ong, *Mathematics*