

MÔ HÌNH TUYỂN TÍNH PHÂN LOẠI HAI LỚP

Lê Thành Sách

✉ Itsach@hcmut.edu.vn

Khoa Khoa học & Kỹ thuật Máy tính
Trường Đại học Bách Khoa - ĐHQG Tp.HCM

Tp.HCM. Ngày 9 tháng 9 năm 2019

Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây
dựng mô hình

Ước lượng tham
số của mô hình

Mục lục

➊ Giới thiệu bài toán

➋ Phương pháp xây dựng mô hình

➌ Ước lượng tham số của mô hình



Hồi quy Logistic

1 Mục lục

Giới thiệu bài toán

Phương pháp xây
dựng mô hình

Ước lượng tham
số của mô hình



Hồi quy Logistic

Mục lục

2 Giới thiệu bài toán

Dữ liệu vào

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Giới thiệu bài toán

Giới thiệu bài toán

Dẫn vào

1 Dữ liệu đầu vào:

$$X = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \cdots & x_{1,(M-1)} \\ 1 & x_{2,1} & x_{2,2} & \cdots & x_{2,(M-1)} \\ 1 & x_{3,1} & x_{3,2} & \cdots & x_{3,(M-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & x_{N,2} & \cdots & x_{N,(M-1)} \end{bmatrix}$$

- Giả sử: dữ liệu gốc đã được rút trích đặc trưng và mỗi điểm dữ liệu là một hàng trong ma trận trên, có M đặc trưng.
- X có kích thước $N \times M$

Giới thiệu bài toán

Dầu vào

② Tập nhãn:

- Bài toán này tập nhãn chỉ có 2 nhãn
- Mỗi nhãn **thường** là tên của lớp; Ví dụ, phân loại trái cây vào “quả táo” và “không phải quả táo”
- \Rightarrow **cần mã hóa nhãn** theo cách dễ xử lý và thuận tiện cho xây dựng mô hình dự báo
- \Rightarrow **Dùng chỉ số**: một nhãn có chỉ số là 0 nhãn kia là 1¹

¹cách quy ước này giúp thuận tiện cho việc xây dựng mô hình ở sau



Hội quy Logistic

Mục lục

Giới thiệu bài toán

4

Dầu vào

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Giới thiệu bài toán

Dầu vào

③ Nhãn của dữ liệu:

- Là một véctơ có N nhãn tương ứng với N điểm dữ liệu.



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

5

Dầu vào

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Giới thiệu bài toán

Dầu vào

Quy ước

- 1 **X** : ma trận dữ liệu; kích thước $N \times M$
- 2 **t** : véctơ chứa nhãn cho dữ liệu; kích thước $N \times 1$. Véctơ này chứa tên nhãn (chuỗi) hay chỉ số.
- 3 **y** : là biểu diễn dạng chỉ số của nhãn t ; kích thước: $N \times 1 \Rightarrow$ có thể dùng **y** và **t** thay thế nhau trong một số trường hợp
- 4 **N** : số điểm dữ liệu
- 5 **M** : số đặc trưng của mỗi điểm dữ liệu
- 6 **\hat{y}** là giá trị dự báo từ mô hình



Hội quy Logistic

Mục lục

Giới thiệu bài toán

6 Dầu vào

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Giới thiệu bài toán

Mục tiêu

Mục tiêu của bài toán là xây dựng mô hình dự đoán từ tập huấn luyện, để sau đó, khi nhận mẫu dữ liệu mới x nó dự đoán nhãn của mẫu này, nghĩa là xuất ra một giá trị thuộc tập nhãn

Bài giảng này hướng dẫn xây dựng mô hình tuyến tính¹ cho bài toán phân lớp.

¹Tuyến tính: chỉ ra rằng đường biên giữa các phân lớp là tuyến tính



Hội quy Logistic

Mục lục

Giới thiệu bài toán

7

Dầu vào

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Hồi quy Logistic

Mục lục

Giới thiệu bài toán

8

Phương pháp xây dựng mô hình

Ý tưởng

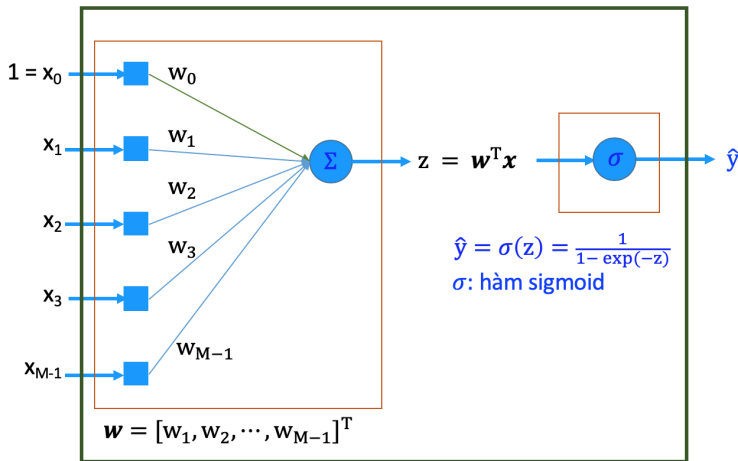
Mô hình dự báo

Ước lượng tham số của mô hình

Phương pháp xây dựng mô hình

Phương pháp xây dựng mô hình

Ý tưởng



Hình 2.1: Mô hình hồi quy logistic 2 lớp

Phương pháp xây dựng mô hình

Ý tưởng



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

10 Ý tưởng

Mô hình dự báo

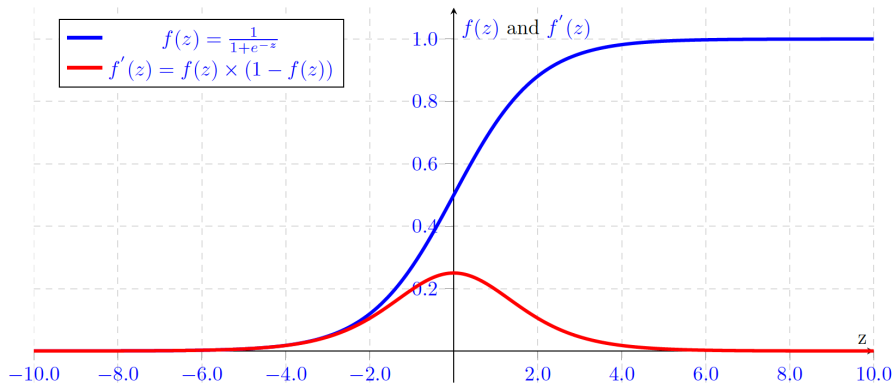
Ước lượng tham số của mô hình

- Hai lớp là C_0 và C_1
- Mô hình nhận vào vectơ \mathbf{x} , tính toán và xuất ra \hat{y} , là xác suất của \mathbf{x} thuộc lớp C_1 , xem Hình 2.1
 - Ngầm hiểu: $1 - \hat{y}$ là xác suất của \mathbf{x} thuộc lớp C_0
- Bên trong
 - 1 Sử dụng một mô hình tuyến tính để dự báo ra điểm $z = \mathbf{w}^T \mathbf{x}^1$
 - 2 Cho z qua hàm **sigmoid** để xuất ra xác suất

¹ \mathbf{w} : tham số của mô hình

Phương pháp xây dựng mô hình

Ý tưởng



Hình 2.2: Hàm sigmoid và đạo hàm

Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

11 Ý tưởng

Mô hình dự báo

Ước lượng tham số của mô hình

Phương pháp xây dựng mô hình

Mô hình dự báo



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ý tưởng

12 Mô hình dự báo

Ước lượng tham số của mô hình

Công thức của mô hình

$$\begin{aligned}\hat{y} &\triangleq p(C_1|\mathbf{x}, \mathbf{w}) \\ &= \sigma(\mathbf{w}^T \mathbf{x}) \\ &= \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}\end{aligned}\tag{2.1}$$

Phương pháp xây dựng mô hình

Mô hình dự báo

- 1 Vào: Dữ liệu x
- 2 Vào: Tham số của mô hình w và ngưỡng λ
- 3 Tính xác suất dự báo \hat{y} theo C.T (2.1)
- 4 Nếu $\hat{y} \geq \lambda$, x thuộc lớp C_1
Ngược lại, x thuộc lớp C_0



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ý tưởng

13

Mô hình dự báo

Ước lượng tham số của mô hình

27

Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây
dựng mô hình

14 Ước lượng tham
số của mô hình

Xây dựng hàm mục tiêu

Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm
bậc 2

Ước lượng tham số của mô hình

Ước lượng tham số của mô hình

Nguyên tắc chung



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

15 **Ước lượng tham số của mô hình**

Xây dựng hàm mục tiêu

Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm bậc 2

Nguyên tắc chung

- Đưa về dạng bài toán tối ưu
 - Sử dụng hàm likelihood
 - Cực đại hóa likelihood cho tập dữ liệu đầu vào
- Giải bài toán tối ưu
 - Không thể phân tích toán học \rightarrow sử dụng phương pháp lặp
 - Gradient Descent
 - Iterative Re-Weighted Least Squares

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

Với một điểm dữ liệu: $\langle \mathbf{x}, y \rangle^1$

Xác suất dự báo được là:

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \hat{y} & \text{if } y = 1 \\ 1 - \hat{y} & \text{if } y = 0 \end{cases}$$

Viết gọn:

$$p(y|\mathbf{x}, \mathbf{w}) = \hat{y}^y (1 - \hat{y})^{1-y} \quad (3.1)$$

¹ Lưu ý: y chỉ nhận giá trị thuộc tập $\{0, 1\}$

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

Với một điểm dữ liệu: $\langle \mathbf{x}_n, y_n \rangle^1$

Xác suất dự báo được là:

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \begin{cases} \hat{y}_n & \text{if } y_n = 1 \\ 1 - \hat{y}_n & \text{if } y_n = 0 \end{cases}$$

Viết gọn:

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \hat{y}_n^{y_n} (1 - \hat{y}_n)^{1-y_n} \quad (3.2)$$

¹ Lưu ý: y_n chỉ nhận giá trị thuộc tập $\{0, 1\}$

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

Xác suất xảy ra N nhãn¹ trong tập huấn luyện:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \hat{y}_n^{y_n} (1 - \hat{y}_n)^{1-y_n} \quad (3.3)$$

Mục tiêu: nguyên lý cực đại hóa hàm hợp lý

- Tìm \mathbf{w} sao cho $p(\mathbf{t}|\mathbf{X}, \mathbf{w})$ đạt cực đại
- Lưu ý: $\hat{y}_n = p(y_n|\mathbf{X}, \mathbf{w})$, phụ thuộc vào \mathbf{w} trong C.T.
(2.1)

¹lấy mẫu theo nguyên tắc i.i.d

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

Sử dụng **negative log-likelihood**:

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &\triangleq -p(\mathbf{t}|\mathbf{X}, \mathbf{w}) \\ &= \sum_{n=1}^N y_n \log \hat{y}_n + (1 - y_n) \log(1 - \hat{y}_n)\end{aligned}\quad (3.4)$$

Mục tiêu: nguyên lý cực đại hóa hàm hợp lý

- Tìm \mathbf{w} sao cho $\mathcal{L}(\mathbf{w})^1$ đạt cực tiểu

¹hàm: **cross-entropy**

Ước lượng tham số của mô hình

Tìm hệ số của mô hình



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

20 Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm bậc 2

Nguyên tắc

- Khó dùng phân tích toán học để giải tìm nghiệm cho bài toán tối ưu có hàm mục tiêu trong C.T (3.4).
- Bài toán cực tiểu hóa $\mathcal{L}(\mathbf{w})$ là tối ưu không ràng buộc, có thể dùng phương pháp lặp.
 - Dựa vào đạo hàm bậc 1: **Gradient Descent**
 - Dựa vào đạo hàm bậc 2: **Newton-Raphson**

Ước lượng tham số của mô hình

Tìm hệ số của mô hình

Cần thiết

Tìm đạo hàm của hàm tổn thất \mathcal{L} so với các tham số w

$$\Delta w \triangleq \frac{\partial \mathcal{L}(w; x, y)}{\partial w}$$



Hội quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

21 Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm bậc 2

Ước lượng tham số của mô hình

Tìm hệ số của mô hình

Phương pháp

Sử dụng **chain-rule**, xem Hình 3.1

$$\frac{\partial \mathcal{L}(w; x, y)}{\partial w} = \frac{\partial \mathcal{L}}{\partial z} \bullet \frac{\partial z}{\partial w}$$

Dấu \bullet là phép “dot”



Hội quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

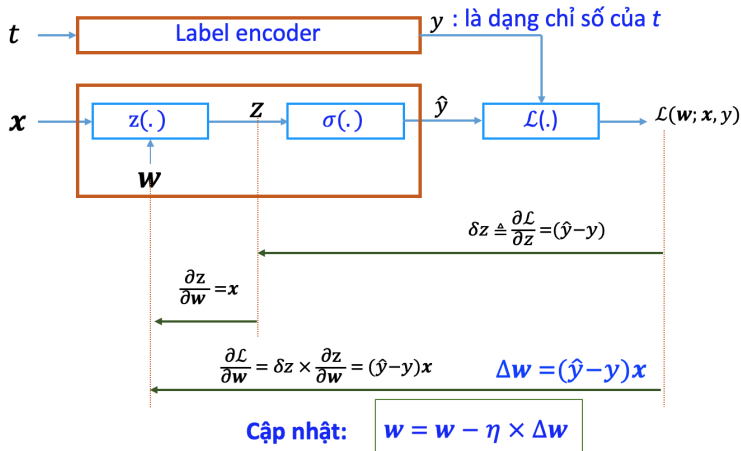
Xây dựng hàm mục tiêu

22 Tìm hệ số của mô hình

Giải thuật lập với đạo hàm bậc 2

Ước lượng tham số của mô hình

Tìm hệ số của mô hình



Hình 3.1: Quy trình tính và cập nhật tham số

Ước lượng tham số của mô hình

Tìm hệ số của mô hình

Đạo hàm của một số hàm trong sơ đồ tính toán, Hình 3.1:

$$\frac{d\mathcal{L}(\mathbf{w}; \mathbf{x}, y)}{d\hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \quad (3.5)$$

$$\frac{d\hat{y}}{dz} = \hat{y}(1 - \hat{y}) \quad (3.6)$$

$$\frac{\partial z}{\partial \mathbf{w}} = \mathbf{x} \quad (3.7)$$

Lưu ý: theo Jacobian, $\frac{\partial z}{\partial \mathbf{w}} = \mathbf{x}^T$; ở đây, phải chuyển vị để phù hợp với kích thước của \mathbf{w} , $M \times 1$



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

24

Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm bậc 2

27

Ước lượng tham số của mô hình

Tìm hệ số của mô hình

Từ các C.T (3.5), (3.6), và (3.7),
Đạo hàm của hàm $\mathcal{L}(\mathbf{w}; \mathbf{x}, y)$ tính trên một điểm dữ liệu
 $\langle \mathbf{x}, y \rangle$ là:

$$\begin{aligned}\Delta \mathbf{w} &= \frac{\partial \mathcal{L}(\mathbf{w}; \mathbf{x}, y)}{\partial \mathbf{w}} \\ &= \frac{d\mathcal{L}(\mathbf{w}; \mathbf{x}, y)}{d\hat{y}} \times \frac{d\hat{y}}{dz} \times \frac{\partial z}{\partial \mathbf{w}} \\ &= (\hat{y} - y)\mathbf{x}\end{aligned}\tag{3.8}$$



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

25 Tìm hệ số của mô hình

Giải thuật lập với đạo hàm bậc 2

Ước lượng tham số của mô hình

Tìm hệ số của mô hình

Đạo hàm của hàm $\mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y})$ tính trên một tập của N điểm dữ liệu $\langle \mathbf{X}, \mathbf{y} \rangle$ là:

$$\begin{aligned}\Delta \mathbf{w} &= \frac{\partial \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y})}{\partial \mathbf{w}} \\ &= \sum_{n=1}^N (\hat{y}_n - y_n) \mathbf{x}_n \\ &= \mathbf{X}^T (\hat{\mathbf{y}} - \mathbf{y})\end{aligned}\tag{3.9}$$



Hồi quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

26 Tìm hệ số của mô hình

Giải thuật lặp với đạo hàm bậc 2

Ước lượng tham số của mô hình

Giải thuật lặp với đạo hàm bậc 2

Đạo hàm bậc 2 của $\mathcal{L}(\mathbf{w})$:

$$\begin{aligned} \mathbf{H} &\triangleq \nabla \nabla \mathcal{L}(\mathbf{w}) \\ &= \sum_{n=1}^N \hat{y}_n (1 - \hat{y}_n) \mathbf{x} \mathbf{x}^T \\ &= \mathbf{X}^T \mathbf{R} \mathbf{X} \end{aligned} \quad (3.10)$$

\mathbf{R} là ma trận đường chéo có phần tử $R_{nn} = \hat{y}_n (1 - \hat{y}_n)$



Hội quy Logistic

Mục lục

Giới thiệu bài toán

Phương pháp xây dựng mô hình

Ước lượng tham số của mô hình

Xây dựng hàm mục tiêu

Tìm hệ số của mô hình

27

Giải thuật lặp với đạo hàm bậc 2

27