

Python for Data Science - Midterm

Nhóm: **Nguyễn Phú Thành**: 18110014 - **Lê Hoàng Đức**: 18110075

13 - 12 - 2020

1 Thuật toán Logistic Regression

Cho bộ dữ liệu $\{\mathbf{x}_n, \mathbf{t}_n\}$, trong đó $\mathbf{t}_n \in \{0, 1\}$ với $n = 1, 2, \dots, N$

Ta gọi C_1 là sự kiện biến ngẫu nhiên $\mathbf{t} = 1$, có xác suất:

$$p(C_1) = \mathbb{P}(\mathbf{t} = 1)$$

và C_2 là sự kiện biến ngẫu nhiên $\mathbf{t} = 0$, có xác suất:

$$p(C_2) = \mathbb{P}(\mathbf{t} = 0)$$

Hơn nữa giả sử vecto đầu vào \mathbf{x} có phân phối chuẩn phụ thuộc vào lớp C_k . Giả sử thêm là các lớp có cùng ma trận hiệp phương sai Σ , tức:

$$p(\mathbf{x}|C_k) = \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma^{-1}(\mathbf{x} - \mu_k) \right\}$$

Khi đó phân phối hậu nghiệm của lớp C_1 dưới dạng:

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(C_1)p(\mathbf{x}|C_1)}{p(\mathbf{x})} \\ &= \frac{p(C_1)p(\mathbf{x}|C_1)}{p(C_1)p(\mathbf{x}|C_1) + p(C_2)p(\mathbf{x}|C_2)} \\ &= \frac{1}{1 + \left(\frac{p(C_2)p(\mathbf{x}|C_2)}{p(C_1)p(\mathbf{x}|C_1)} \right)} \\ &= \frac{1}{1 + e^{-a}} = \sigma(a) \end{aligned}$$

trong đó:

$$\begin{aligned} a &= \ln \left(\frac{p(C_1)p(\mathbf{x}|C_1)}{p(C_2)p(\mathbf{x}|C_2)} \right) \\ &= \ln \left(\frac{p(C_1)}{p(C_2)} \right) - \frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2) \\ &= \ln \left(\frac{p(C_1)}{p(C_2)} \right) + \frac{1}{2}(\mu_2^T \Sigma^{-1} \mu_2 - \mu_1^T \Sigma^{-1} \mu_1) + (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} \\ &= \mathbf{w}^T \mathbf{x} + w_0 \end{aligned}$$

với

$$\begin{aligned} \mathbf{w} &= \Sigma^{-1}(\mu_1 - \mu_2) \\ w_0 &= \ln \left(\frac{p(C_1)}{p(C_2)} \right) + \frac{1}{2}(\mu_2^T \Sigma^{-1} \mu_2 - \mu_1^T \Sigma^{-1} \mu_1) \end{aligned}$$

Đặt $\tilde{\mathbf{w}} = (w_0, \mathbf{w})^T$, $\tilde{\mathbf{x}} = (1, \mathbf{x})^T$, $\mathbf{y}(\mathbf{x}) = \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$

Vậy:

$$\begin{aligned} p(C_1|\mathbf{x}) &= \mathbf{y}(\mathbf{x}) \\ p(C_2|\mathbf{x}) &= 1 - \mathbf{y}(\mathbf{x}) \end{aligned}$$

với $\sigma(x) = \frac{1}{1 + e^{-x}}$ là hàm sigmoid

Ta có hàm hợp lý:

$$\begin{aligned} p(\mathbf{t}|\mathbf{x}, \tilde{\mathbf{w}}) &= \prod_{n=1}^N p(t_n|\mathbf{x}_n, \tilde{\mathbf{w}}) \\ &= \prod_{n=1}^N \mathbf{y}_n^{t_n} (1 - \mathbf{y}_n)^{1-t_n} \end{aligned}$$

trong đó: $\mathbf{y}_n = y(\mathbf{x}_n) = \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)$ và $\mathbf{t} = (t_1, t_2, \dots, t_N)^T$
Lấy đối của hàm ln của hàm hợp lý, ta được:

$$E_D(\tilde{\mathbf{w}}) = -\ln p(\mathbf{t}|\mathbf{x}, \tilde{\mathbf{w}}) = -\sum_{n=1}^N [t_n \ln \mathbf{y}_n + (1 - t_n) \ln(1 - \mathbf{y}_n)]$$

Ta gọi hàm E_D là hàm sai số cross-entropy

Để tìm ước lượng hợp lý cực đại, ta tìm giá trị nhỏ nhất của hàm sai số cross - entropy này

Vì hàm sai số E_D là hàm lồi nên để tìm giá trị nhỏ nhất, ta đưa bài toán về việc tìm nghiệm của phương trình $\nabla E_D(\tilde{\mathbf{w}}) = 0$

Ta có:

$$\begin{aligned} d\mathbf{y}_n &= d\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \\ &= \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)(1 - \sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n))d(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \\ &= \mathbf{y}_n(1 - \mathbf{y}_n)\tilde{\mathbf{x}}_n d\tilde{\mathbf{w}} \end{aligned}$$

Khi đó:

$$\begin{aligned} dE_D &= -\sum_{n=1}^N t_n \frac{1}{\mathbf{y}_n} d\mathbf{y}_n + (1 - t_n) \frac{-1}{1 - \mathbf{y}_n} d\mathbf{y}_n \\ &= -\sum_{n=1}^N t_n \frac{1}{\mathbf{y}_n} \mathbf{y}_n(1 - \mathbf{y}_n)\tilde{\mathbf{x}}_n d\tilde{\mathbf{w}} + (1 - t_n) \frac{-1}{1 - \mathbf{y}_n} \mathbf{y}_n(1 - \mathbf{y}_n)\tilde{\mathbf{x}}_n d\tilde{\mathbf{w}} \\ &= -\sum_{n=1}^N [t_n(1 - \mathbf{y}_n)\tilde{\mathbf{x}}_n - (1 - t_n)\mathbf{y}_n\tilde{\mathbf{x}}_n] d\tilde{\mathbf{w}} \\ &= \sum_{n=1}^N (\mathbf{y}_n - t_n)\tilde{\mathbf{x}}_n d\tilde{\mathbf{w}} \end{aligned}$$

Suy ra:

$$\nabla E_D(\tilde{\mathbf{w}}) = \sum_{n=1}^N (\mathbf{y}_n - t_n)\tilde{\mathbf{x}}_n = \sum_{n=1}^N [\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) - t_n] \tilde{\mathbf{x}}_n$$

Sau đây, ta sẽ trình bày một trong nhiều cách để giải phương trình $\frac{\partial E_D}{\partial \tilde{\mathbf{w}}} = 0$

2 Iterative reweighted least squares (IRLS) method

Cách cập nhật $\tilde{\mathbf{w}}$ để hàm E_D đạt giá trị nhỏ nhất, được sử dụng theo phương pháp Newton có dạng:

$$\tilde{\mathbf{w}}^{(\tau+1)} = \tilde{\mathbf{w}}^{(\tau)} - \eta \nabla^2 E_D(\tilde{\mathbf{w}}^{(\tau)})^{-1} \nabla E_D(\tilde{\mathbf{w}}^{(\tau)})$$

trong đó η là tốc độ học
Ta có:

$$\begin{aligned}\nabla E_D(\tilde{\mathbf{w}}) &= \sum_{n=1}^N (\mathbf{y}_n - \mathbf{t}_n) \tilde{\mathbf{x}}_n = \mathbf{X}^T (\mathbf{y} - \mathbf{t}) \\ \nabla^2 E_D(\tilde{\mathbf{w}}) &= \nabla \sum_{n=1}^N (\mathbf{y}_n - \mathbf{t}_n) \tilde{\mathbf{x}}_n \\ &= \sum_{n=1}^N \mathbf{y}_n (1 - \mathbf{y}_n) \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^T = \mathbf{X}^T \mathbf{R} \mathbf{X}\end{aligned}$$

trong đó \mathbf{R} là ma trận đường chéo có hệ số trên đường chéo chính là:

$$R_{nn} = \mathbf{y}_n (1 - \mathbf{y}_n)$$

Thay vào ta được:

$$\begin{aligned}\tilde{\mathbf{w}}^{(\text{new})} &= \tilde{\mathbf{w}}^{(\text{old})} - \eta (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{t}) \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{R} \mathbf{X} \tilde{\mathbf{w}}^{(\text{old})} - \eta \mathbf{X}^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{R} \mathbf{z}\end{aligned}$$

trong đó: $\mathbf{z} = \mathbf{X} \tilde{\mathbf{w}}^{(\text{old})} - \eta \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$

3 Stochastic gradient descent (SGD) method

Tại mỗi vòng lặp, $\tilde{\mathbf{w}}$ được cập nhật dựa trên một điểm dữ liệu ngẫu nhiên. Hàm mất mát E_D của hồi quy logistic với một điểm dữ liệu $(\mathbf{x}_n, \mathbf{t}_n)$ và gradient của nó ta thu được từ trình bày phía trên

$$\nabla_{\tilde{\mathbf{w}}} E_D = (\mathbf{y}_n - \mathbf{t}_n) \tilde{\mathbf{x}}_n = (\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) - \mathbf{t}_n) \tilde{\mathbf{x}}_n$$

Từ đó, công thức cập nhật nghiệm cho hồi quy logistic sử dụng SGD là

$$\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} - \eta (\mathbf{y}_n - \mathbf{t}_n) \tilde{\mathbf{x}}_n = \tilde{\mathbf{w}} - \eta (\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) - \mathbf{t}_n) \tilde{\mathbf{x}}_n$$

với η là tốc độ học.

4 Hồi quy logistic với chính quy hóa

Một trong các kỹ thuật phổ biến giúp vectơ trọng số $\tilde{\mathbf{w}}$ không "quá lớn" (dưới một chuẩn nào đó) là sử dụng kỹ thuật chính quy hóa (Regulization).

Đây là một kỹ thuật kiểm soát, trong đó một đại lượng tỉ lệ với bình phương chuẩn l_2 của vector trọng số $\tilde{\mathbf{w}}$ được cộng vào hàm mất mát để kiểm soát độ lớn của các hệ số.

Khi đó tối ưu hàm tổng mất mát được viết dưới dạng

$$E(\tilde{\mathbf{w}}) = E_D(\tilde{\mathbf{w}}) + \lambda E_{\tilde{\mathbf{W}}}(\tilde{\mathbf{w}})$$

trong đó, λ là hệ số chuẩn hóa để kiểm soát mối tương quan giữa hàm mất mát phụ thuộc vào dữ liệu huấn luyện $E_D(\tilde{\mathbf{w}})$ và số hạng kiểm soát $E_{\tilde{\mathbf{W}}}(\tilde{\mathbf{w}})$.

Hai hàm kiểm soát phổ biến nhất là l_1 norm và l_2 norm. Ví dụ, khi chọn $E_{\tilde{\mathbf{W}}}(\tilde{\mathbf{w}}) = \frac{1}{2} \|\tilde{\mathbf{w}}\|_2^2$

$$E_{\tilde{\mathbf{W}}}(\tilde{\mathbf{w}}) = \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}}$$

thì hàm mất mát trở thành

$$E(\tilde{\mathbf{w}}) = - \sum_{n=1}^N [\mathbf{t}_n \ln \mathbf{y}_n + (1 - \mathbf{t}_n) \ln(1 - \mathbf{y}_n)] + \frac{\lambda}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}}$$

Tính toán tương tự trên ta sẽ được:

$$\begin{aligned} \nabla E_D(\tilde{\mathbf{w}}) &= \sum_{n=1}^N (\mathbf{y}_n - \mathbf{t}_n) \tilde{\mathbf{x}}_n + \lambda \tilde{\mathbf{w}} = \sum_{n=1}^N [\sigma(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) - \mathbf{t}_n] \tilde{\mathbf{x}}_n + \lambda \tilde{\mathbf{w}} \\ \nabla^2 E_D(\tilde{\mathbf{w}}) &= \mathbf{X}^T \mathbf{R} \mathbf{X} + \lambda \end{aligned}$$

Thay các giá trị này vào các công thức trên ta được cách cập nhật cho tham số $\tilde{\mathbf{w}}$