

Họ và tên: Đặng Anh Tiến
MSSV: 20520800
Lớp: CS431.N11

Bài tập LT2: Giải thích hàm loss function của Logistic Regression

Trong bài toán Binary Classification, miền giá trị của y_{pred} (kết quả dự đoán) nằm trong khoảng $(0, 1)$; y_{true} (kết quả thực tế) là giá trị 0 và 1.

- Trong trường hợp kết quả dự đoán gần với kết quả thực tế ($y_{\text{true}} == y_{\text{pred}}$):
 $y_{\text{true}} = 0, y_{\text{pred}} \sim 0$

$$\text{MSE} = (0 - 0)^2 = 0$$

$$\text{BCE} = -(0 * \log(0) + \log(1)) = 0$$

=> Cho thấy khi mô hình dự đoán đúng, kết quả của hai hàm loss không có quá nhiều sự khác biệt.

- Trong trường hợp kết quả dự đoán khác với kết quả thực tế ($y_{\text{true}} != y_{\text{pred}}$):
 $y_{\text{true}} = 0, y_{\text{pred}} \sim 1$

$$\text{MSE} = (0 - 1)^2 \sim 1$$

$$\text{BCE} = -(0 * \log(1) + \log(0)) = -(-\infty) = +\infty$$

=> BCE trả về $+\infty$ cao hơn rất nhiều so với MSE.

Có thể thấy trong bài toán Binary Classification, ta nên dùng Binary Cross Entropy (BCE) thay vì Mean Squared Error (MSE) bởi vì các lý do sau:

- Giá trị của y_{true} và y_{pred} trong khoảng $[0, 1]$, vì vậy MSE luôn trả về giá trị trong khoảng $[0, 1]$. Khi model dự đoán sai, BCE, như đã chứng minh ở trên, trả về kết quả lớn hơn rất nhiều ($+\infty$) so với MSE. Điều này khiến BCE phạt hàm Loss nặng hơn.
- Giá trị hàm Loss lớn hơn, gradient lớn hơn, vì vậy trong Gradient Descent, tối ưu các parameter sẽ diễn ra nhanh hơn. MSE có giá trị trong $[0, 1]$, các gradient sẽ có giá trị $[0, 1]$, Gradient Descent sẽ tối ưu vô cùng chậm. Trong khi đó, BCE cho kết quả Loss lớn, gradient sẽ lớn, điều này giúp Gradient Descent tối ưu nhanh hơn nhiều so với MSE.