

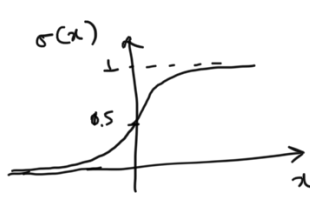
Logistic regression (classification)

Example: $p(s|x)$, $x = (x_1, x_2, x_3)$, x_1 : age, x_2 : m/f, x_3 : chok

The simplest approach would assume some linear model:

(linear regression) $w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 = w^T x$, $x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$, $w^T = [w_0, w_1, w_2, w_3]$

... but this is not a probability! We can fix this by "warping" $w^T x$ with a sigmoid function:


$$\left. \begin{aligned} \sigma(x) &= \frac{1}{1 + e^{-x}} \\ \text{logistic sigmoid} \end{aligned} \right\} p(y|x) = \sigma(w^T x)$$

Formal definition:

Setup: Given $\mathcal{D} := \{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \{0, 1\}$

Model: $y_i \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(\underbrace{\sigma(w^T x_i)}_{a_i})$, $\theta := \{w_0, w_1, \dots, w_d\}$

Pros: • Interpretable (= model parameters are meaningful)

• Reveal which variables are more influential,

• Small number of trainable parameters ($d+1$)

↳ simple model that is "statistically easy to train"

• Computationally efficient ways to estimate w .

• Extension to multi-class is straightforward.

• Forms the foundation for more complex models (GLM, NNs)

Cons:

- Being a simple model, its performance is inferior to more complex models.

Maximum Likelihood Estimation:

$$\underline{w}_{MLE} = \arg \max_w p(y|X, w)$$

$$a_i := \sigma(w^T x_i)$$

$$p(y_1, \dots, y_n | x_1, \dots, x_n, w_0, w_1, \dots, w_d) \stackrel{i.i.d.}{=} \prod_{i=1}^n p(y_i | x_i, w)$$

$$= \prod_{i=1}^n a_i^{y_i} (1 - a_i)^{1 - y_i}$$

$$\underline{L(w)} := -\log p(y|X, w) = - \sum_{i=1}^n y_i \log a_i + (1 - y_i) \log (1 - a_i) \quad \left. \vphantom{\sum_{i=1}^n} \right\} \text{Binary cross-entropy loss}$$

• Recall that $a_i = \sigma(w^T x_i)$

$$\rightarrow \log a_i = -\log(1 + e^{-w^T x_i})$$

$$\rightarrow \log(1 - a_i) = -w^T x_i - \log(1 + e^{-w^T x_i})$$

$$\rightarrow \frac{\partial}{\partial w_i} \log a_i = x_i (1 - a_i)$$

$$\rightarrow \frac{\partial}{\partial w_i} \log(1 - a_i) = -a_i x_i$$

$$\begin{aligned} \Rightarrow \nabla_{w_j} L(w) &= - \sum_{i=1}^n y_i x_{ij} (1 - a_i) - (1 - y_i) x_{ij} a_i \\ &= \sum_{i=1}^n (a_i - y_i) x_{ij} \end{aligned}$$

$$\boxed{\nabla_w L(w) = X^T (a - y)} \quad , \quad X = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix}_{n \times (d+1)}, \quad a = \begin{bmatrix} \sigma(w^T x_1) \\ \vdots \\ \sigma(w^T x_n) \end{bmatrix}$$
$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Lyns

$$\Downarrow$$

$$\boxed{\nabla_w^2 \mathcal{L}(w) = X^T A X}, \quad A := \begin{bmatrix} a_1(1-a_1) & & 0 \\ & \ddots & \\ 0 & & a_n(1-a_n) \end{bmatrix}_{n \times n}$$

↪ symmetric PSD

Newton's method: (Iterative re-weighted least squares)

$$w_{t+1} = w_t - \eta H_t^{-1} g_t, \quad H_t := X^T A_t X, \quad g_t := X^T (a_t - y)$$

$$\Rightarrow w_{t+1} = w_t - (X^T A_t X)^{-1} X^T (a_t - y)$$

Rewrite as:

$$w_{t+1} = (X^T A_t X)^{-1} X^T A_t \left[\overbrace{X w_t - A_t^{-1} (a_t - y)}^{z_t} \right]$$

$$\Rightarrow w_{t+1} = \boxed{(X^T A_t X)^{-1} X^T A_t z_t} \quad \left(\begin{array}{l} \text{Newton for logistic} \\ \swarrow \end{array} \right)$$

Recall
Linear regression: $\underline{w_{MLE}} = \underbrace{(X^T X)^{-1}}_{\leftarrow} X^T y = (X^T A X)^{-1} X^T A y$,
where $\underline{A := I}$.

Multi-class logistic regression:

$$\text{Model: } \underbrace{p(y=c|X, w)}_{\substack{c=1, \dots, q \\ \hookrightarrow \# \text{ of classes}}} = \frac{\exp(w_c^T x)}{\sum_{c'=1}^q \exp(w_{c'}^T x)} \quad \left. \begin{array}{l} \text{soft-max is} \\ \text{a generalization} \\ \text{of the logistic} \\ \text{sigmoid for} \\ \text{multiple classes.} \end{array} \right\}$$

where w_c is the c -th column of W
($d \times 1$) \times q ,

and y is a "one-hot" encoding vector: $\begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$
... i.e. $y_i = c_i$

$n \times c$

$$\underline{\underline{y_{ic}}} = \mathbb{1}_{\{y_i=c\}} = \begin{cases} 1, & \text{if } y_i = c \\ 0, & \text{if } y_i \neq c \end{cases} \quad \left[\begin{matrix} \text{ } \\ \text{ } \\ \text{ } \\ \text{ } \end{matrix} \right]_c$$

$$p(y|x, w) = \prod_{i=1}^n \prod_{c=1}^d p(y_i = c | x_i, w)$$

$$\Rightarrow -\log p(y|x, w) = \sum_{i=1}^n \left[\sum_{c=1}^d y_{ic} w_c^T x_i \right] - \log \left(\sum_{c=1}^d \exp(w_c^T x_i) \right)$$

→ multi-class cross entropy.