

# Logistic Regression

Classification:  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}, x \in \mathbb{R}^D, y \in \{0, 1\}$

Linear classifier in higher dimension:

Hyperplane:  $H = \{x : w^T x = 0\}, H^+ = \{x : w^T x > 0\}, H^- = \{x : w^T x < 0\}$

Prediction using a decision boundary:

$$h(x) = \begin{cases} 1 & w^T x \geq 0 \\ 0 & w^T x < 0 \end{cases}$$

Estimating probabilities — logistic (sigmoid) function

$$z(x) = w^T x \rightarrow (-\infty, +\infty)$$

$$\sigma(z(x)) = \frac{1}{1+e^{-z(x)}} \rightarrow (0, 1)$$

How can we find best hyperplane  $w$ ?

Data  $\rightarrow$  Estimation

T, H, T, H, T  $\rightarrow$  to predict  $\theta$  (probability of head), find  $\theta$  that maximizes  $(1 - \theta)\theta(1 - \theta)\theta(1 - \theta)$

Likelihood function  $L(\theta) = p(D|\theta) = \theta^{N_H} (1 - \theta)^{N_T}, l(\theta) = \ln(L(\theta))$

**Maximum Likelihood Estimation (MLE):** maximize  $l(\theta)$

Extend this to conditional likelihood:

$$p(y = 1|x; w) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}$$

$$p(y = 0|x; w) = 1 - \sigma(w^T x) = 1 - \frac{1}{1+e^{-w^T x}}$$

$$L(w) = \prod_{i=1}^N \sigma(w^T x^{(i)})^{y^{(i)}} (1 - \sigma(w^T x^{(i)}))^{1-y^{(i)}}$$

$$l(w) = \sum_{i=1}^N [y^{(i)} \ln \sigma(w^T x^{(i)}) + (1 - y^{(i)}) \ln(1 - \sigma(w^T x^{(i)}))]$$

Example:  $y^{(i)} = 0, \sigma(w^T x^{(i)}) = 0 \rightarrow \ln(1) = 0$

$$y^{(i)} = 0, \sigma(w^T x^{(i)}) = 0.99 \rightarrow \ln(0.01) = -4.61$$

## Gradient Ascent

We want to maximize  $\frac{1}{N}l(w)$

$$w^* = \arg \max_w \left( \frac{1}{N} \sum_{i=1}^N [y^{(i)} \ln \sigma(w^T x^{(i)}) + (1 - y^{(i)}) \ln(1 - \sigma(w^T x^{(i)}))] \right)$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)(1 - \sigma(z))$$

$$\frac{\partial \sigma(w^T x)}{\partial w_j} = \frac{\partial \sigma(w_0 x_0 + \dots + w_d x_d)}{\partial w_j} = x_j$$

$$\frac{\partial \sigma(z)}{\partial w_j} = \frac{\partial \sigma(w^T x)}{\partial w^T x} \cdot \frac{\partial w^T x}{\partial w_j} = \sigma(w^T x)(1 - \sigma(w^T x))x_j$$

$$\frac{\partial}{\partial w_j} l(w) = \sum_{i=1}^N (y_i - \sigma(w^T x^{(i)})) x_j^{(i)}$$

If  $y_i \approx \sigma(w^T x^{(i)})$ , almost no change!

If  $y_i - \sigma(w^T x^{(i)}) \approx \pm 1$ , approx  $\frac{\alpha}{N}$  times the  $j^{th}$  feature!

for i = 1 to num\_iter:

$$temp0 = w_0 + \frac{\alpha}{N} \frac{\partial l(w)}{\partial w_0} = w_0 + \frac{\alpha}{N} \sum_{i=1}^N (y_i - \sigma(w^T x^{(i)})) x_0^{(i)} \quad // + \text{ since ascent}$$

$$temp1 = w_1 + \frac{\alpha}{N} \frac{\partial l(w)}{\partial w_1} = w_1 + \frac{\alpha}{N} \sum_{i=1}^N (y_i - \sigma(w^T x^{(i)})) x_1^{(i)}$$

...

$$tempd = w_d + \frac{\alpha}{N} \frac{\partial l(w)}{\partial w_d} = w_d + \frac{\alpha}{N} \sum_{i=1}^N (y_i - \sigma(w^T x^{(i)})) x_d^{(i)}$$

$$w_0 = temp0$$

$$w_1 = temp1$$

...

$$w_d = tempd$$

## Vector Implementation

for i = 1 to num\_iter:

$$w = w + \frac{\alpha}{N} X^T (y - \sigma(Xw))$$

## Evaluating errors: Precision and Recall

Two types of error:

1. False positive — predict positive (has cancer), but false
2. False negative — predict negative, but false

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{F1 Score} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \in [0, 1]$$

## Regularization of Logistic Regression

$$l_{lasso}(w) = \frac{1}{N} l(w) - \lambda (\|w_{1:d}\|_1) \quad // \text{ since ascent}$$

$$l_{ridge}(w) = \frac{1}{N} l(w) - \lambda (\|w_{1:d}\|_2^2)$$

## Multiple Classes ( $C_1, \dots, C_k$ )

One-versus-One approach:  $\frac{K(K-1)}{2}$  binary classification problems

classify into  $C_i, C_j$ ; predict the class that wins "majority of votes", confidence scores to resolve ties

One-versus-All approach:  $K$  binary classification problems

classify into  $C_i$  and  $all - C_i$ ; predict the class that has largest confidence score

Could our algorithm directly estimate the probability of label belonging to each of the classes?

(i.e. don't resort to a binary classification problem)

We will predict  $K$  different probabilities:  $y^{(i)} = [y_1^{(i)}, \dots, y_K^{(i)}]^T$

Logistic regression:  $p(y = 1|x; w) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}} = \frac{e^{w^T x}}{e^{w^T x} + 1}$

Soft-max:  $p(y = j|x; w) = \frac{e^{w_j^T x}}{\sum_{j=1}^K e^{w_j^T x}}$