

LOGISTIC REGRESSION - I

Classification : Spam vs. Normal Email
 Malignant vs. Healthy Tissue
 Fake vs. Real News.
 Human vs. Animal image.

One of the most popular ML algos.

$y_i \in \{0, 1\}$ \rightarrow True/Yes
 \rightarrow False/No

No classifier is perfect.

So, we are looking for a probability estimate for $p(y_i/x_i; \theta)$

$\hat{y}_i = h(x_i; \theta)$ \rightarrow hypothesis
 & $0 \leq h \leq 1 \quad \forall x_i, \theta$

Decision boundary:
 (Linear) \rightarrow vector
 or HYPERPLANE $\theta_0 + \theta^T x_i > 0 \Rightarrow$

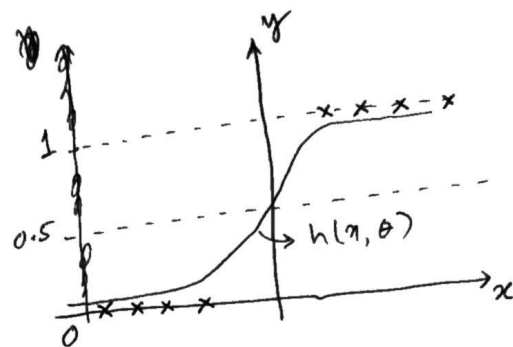
$\theta_0 + \theta^T x_i \leq 0 \Rightarrow$

$\left. \begin{array}{l} h(x_i; \theta) > 0.5 \\ h(x_i; \theta) \leq 0.5 \end{array} \right\} \Rightarrow h(x_i; \theta) = h(\theta^T x_i + \theta_0)$

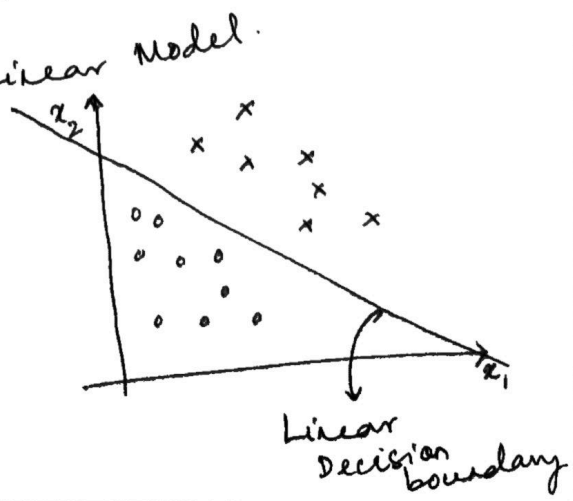
Sigmoid/Logistic function:
 or Logit

many other options available but this one is easier to handle & interpret.

$\log \frac{h(x_i; \theta)}{1 - h(x_i; \theta)}$ odds $= \theta_0 + \theta^T x_i$: Linear Model.

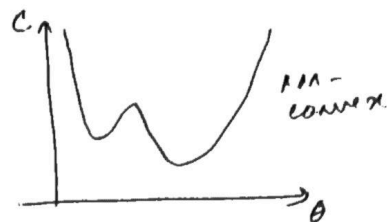


Clearly y vs. x is NOT a straight line. So, makes no sense to ~~fit a~~ use linear regression



$$\text{LSE} : C_{\text{LSE}}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

But this not convex.



for a function, $f(\theta)$, to be convex.

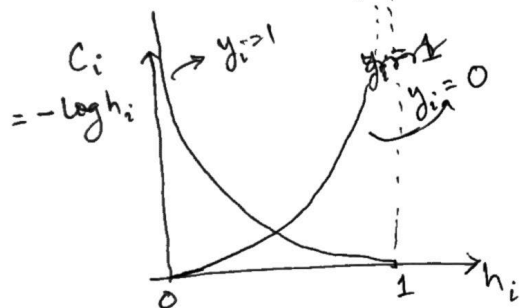
we must have $\frac{\partial^2 f}{\partial \theta^2} > 0 \quad \forall \theta$.

for binary classification,

$$\theta_{\text{MLE}} = \underset{\theta}{\text{argmin}} \quad \frac{1}{N} \sum_{i=1}^N [-y_i \log h_i - (1-y_i) \log (1-h_i)]$$

$$\therefore C(\theta) = \frac{1}{N} \sum_{i=1}^N [-y_i \log h_i - (1-y_i) \log (1-h_i)] \quad ; y_i \in \{0, 1\}$$

↳ Binary Cross Entropy.



if $y_i=1$, $C_i(\theta) = -\log h_i$
 if $y_i=0$, $C_i(\theta) = -\log (1-h_i)$

Gradient Descent:

$$\theta_j \rightarrow \theta_j - \eta \frac{\partial C}{\partial \theta_j}$$

↳ Learning Rate

LOGISTIC REGRESSION - II

MULTI-CLASS CLASSIFICATION

$S = \{\text{spam, Important, Not Important}\}$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

ONE-HOT
ENCODING

$$\hat{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \hat{y}^{(3)} \end{bmatrix}$$

$$\begin{aligned} \hat{y}^{(1)} &= P[y^{(1)} = 1] \\ \hat{y}^{(2)} &= P[y^{(2)} = 1] \\ \hat{y}^{(3)} &= P[y^{(3)} = 1] \end{aligned}$$

$$0 \leq \hat{y}^{(k)} \leq 1$$

$$\sum_k \hat{y}^{(k)} = 1$$

SOFTMAX function:

$$\hat{y}^{(1)} =$$

$$\hat{y}^{(2)} =$$

$$\frac{z^{(2)}}{z^{(1)} + z^{(2)} + z^{(3)}}$$

$$\hat{y}^{(3)} = \frac{z^{(3)}}{z^{(1)} + z^{(2)} + z^{(3)}}$$

$$z^{(k)} = e^{x^T \theta^{(k)} + \theta_0^{(k)}}$$

COST / LOSS FUNCTION

for 2-classes:
$$C = \frac{1}{N} \sum_{i=1}^N [-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)]$$

$$= \frac{1}{N} \sum_{i=1}^N [-y_i^{(1)} \log \hat{y}_i^{(1)} - y_i^{(2)} \log \hat{y}_i^{(2)}]$$

for K-classes:
$$C = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K -y_i^{(k)} \log \hat{y}_i^{(k)}$$

CATEGORICAL
CROSS ENTROPY

LOGISTIC REGRESSION - I (cont'd.)

Classification: $y_i \in \{0, 1\}$

$$\hat{y}_i = h(x_i, \theta) = \frac{1}{1 + e^{-(\theta^T x_i + \theta_0)}} \rightarrow \text{SIGMOID, LOGISTIC, LOGIT f.t.}$$

$$C(\theta) = \frac{1}{N} \sum_{i=1}^N [-y_i \log h_i - (1 - y_i) \log (1 - h_i)] \rightarrow \text{BINARY CROSS ENTROPY.}$$

Gradient Descent: $\theta_j \rightarrow \theta_j - \eta \frac{\partial C}{\partial \theta_j}$ Learning Rate.

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N \left[-\frac{y_i}{h_i} \frac{\partial h_i}{\partial \theta_j} + \frac{1 - y_i}{1 - h_i} \frac{\partial h_i}{\partial \theta_j} \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{(h_i - y_i)}{h_i(1 - h_i)} \frac{\partial h_i}{\partial \theta_j}$$

$$= \frac{1}{N} \sum_{i=1}^N (h_i - y_i) x_{ij}$$

Form Exactly same as Linear Regression just with a different expression for h_i

$$\frac{\partial h_i}{\partial \theta_j} = \frac{-1}{(1 + e^{-\theta^T x_i - \theta_0})^2} \cdot \frac{\partial e^{-\theta^T x_i - \theta_0}}{\partial \theta_j}$$

$$= \frac{e^{-\theta^T x_i - \theta_0}}{(1 + e^{-\theta^T x_i - \theta_0})^2} \cdot \frac{\partial e^{-\theta^T x_i - \theta_0}}{\partial \theta_j}$$

$$= \frac{1}{1 + e^{-\theta^T x_i - \theta_0}} \cdot \frac{e^{-\theta^T x_i - \theta_0}}{1 + e^{-\theta^T x_i - \theta_0}} \cdot x_{ij}$$

$$= h_i(1 - h_i) x_{ij}$$

The features x_{ij} used in Logit can actually be nonlinear f.t. (eg. polynomial) of other features thereby allowing the decision boundary to be effectively nonlinear.

REGULARISATION

$$C(\theta) = \frac{1}{N} \sum_{i=1}^N [-y_i \log h_i - (1 - y_i) \log (1 - h_i)] + \frac{\lambda}{2N} \|\theta\|_2^2$$

as $N \rightarrow \infty$, this term drops out.

$$\therefore \frac{\partial C}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N (h_i - y_i) x_{ij} + \frac{\lambda}{N} \theta_j$$

MULTI-CLASS LOGISTIC REGRESSION

Email: Spam, Important, Not-Important
Face Identification

ONE vs. ALL [or one vs. Rest]

Let $y_i \in \{S, I, NI\} \equiv K$ -categories

Divide into three problems

$$h^{(1)}(x_i, \theta) = \frac{1}{S} \text{ vs. } \frac{0}{\{I, NI\}}$$

$$h^{(2)}(x_i, \theta) = I \text{ vs. } \{S, NI\}$$

$$h^{(3)}(x_i, \theta) = NI \text{ vs. } \{S, I\}$$

K-binary classifications

$$h(x_i, \theta) = \max \{h^{(1)}, h^{(2)}, h^{(3)}\}$$

ONE vs. ONE

S vs. I

I vs. NI

S vs. NI

K_{C2} binary classifications.

Pick category which has most classifications.

These methods do not give ~~good~~ probability estimates
since $h^{(1)} + h^{(2)} + h^{(3)} \neq 1$

One vs. One is also ~~hard~~ computationally expensive.
Better to have a method which can learn multi-class parameter at one go.