# Logistic Regression

Logistic regression is used when the target variable is categorical.

ex:
    to predict email is spam (1) or ham (1).

    to predict tumor is malignant (1) or benign (0).

If we use linear regression, there is a need for setting up a threshold based on which the classification could be done.

It can be observed that having a fixed threshold is not a good way to make predictions in real time.

Linear regression is unbounded and this brings logistic regression into picture, their value ranges from 0 to 1.

## Model structure:

$$\text{output} = 0 \text{ or } 1$$

$$\text{Hypothesis}: \quad z = w \cdot x + \beta$$

$$h\theta(x) = \text{sigmoid}(z)$$

sigmoid $\Rightarrow \dfrac{1}{e^{-x} + 1}$



if $z \to \infty$, $\hat{y} = 1 \to c_1$

$z \to -\infty$, $\hat{y} = 0 \to c_2$

predicted value.

we can see that we first fit into linear regression model, which is acted upon by an activation function / logistic function predicting the target categorical dependent variable.

Types of Logistic regression:

① Binary Logistic regression — only 2 possible outcomes.

② Multinomial Logistic regression — 3 or more categories without ordering. [vegan, non-vegan ...]

③ Ordinal Logistic " — 3 or more categories with ordering [move rating 1-10].

# Decision Boundary:

To predict which class, a data belongs to, a threshold can be set. Based upon this threshold, the obtained estimated probability is classified into classes.

Note: Decision boundary can be linear | non-linear. The polynomial order can be increased to get complex decision boundary.

cost function:

$$cost(h_\theta(x), y_{actual}) = -\log(h_\theta(x)) \text{ if } y=1$$

(Q) why can't we use the cost function used for Linear regression?

$$= -\log(1 - h_\theta(x)) \text{ if } y=0$$

Linear regression uses MSE as its cost function. If this is used for Logistic regression, then it will be a non-convex function.

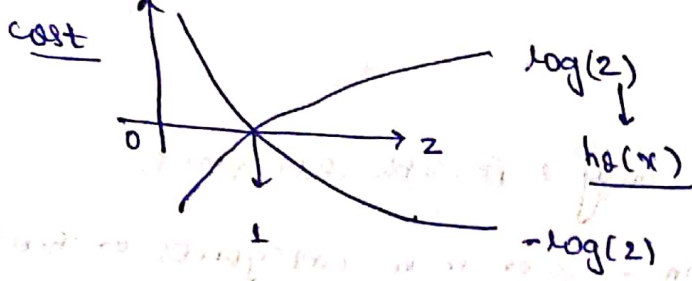Gradient descent will converge into global min. only if the function is convex.
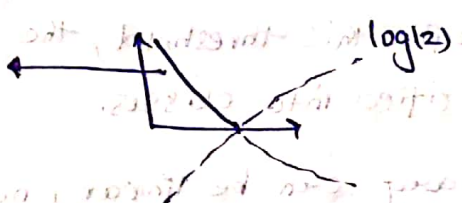


Convex            non-convex

cost



log(2)

$h_\theta(x)$

$-log(2)$

$cost(h_\theta(x), y) = -log(h_\theta(x)) \Rightarrow$ if $y=1$

$= -log(1 - h_\theta(x)) \Rightarrow$ if $y = 0$

if $y = 1:$  $cost(h_\theta(x), y) = -log(h_\theta(x))$

cost



$h_\theta(x)$

$log(2)$

if cost $= 0 \Rightarrow y = 1 \Rightarrow h_\theta(x) = 1$

cost $= \infty \Rightarrow h_\theta(x) = 0$

cost



$h_\theta(x)$

$-log(1-2)$

if cost $= 0 \Rightarrow y = 0 \Rightarrow h_\theta(x) = 0$

cost $= \infty \Rightarrow y \to 1 \Rightarrow h_\theta(x) = 1$

if $h_\theta(x) = 1$, it is similar to predicting $P(y=0 \mid x; \theta) = 0$

simplified cost function:

$$cost(h_\theta(x), y) = -y \cdot log(h_\theta(x)) - (1-y) log(1 - h_\theta(x))$$

notes

we use negation because when we train, we need to maximize the prob. by min. the loss function.

for classification into more than 2 classes, softmax activation function can be used.

---

∎ cost vs Loss

The terms cost and loss functions refer to the same meaning.

Loss function mainly applies for a single training set as compared to the cost function which deals with a penalty for a number of training sets or a complete batch. Also known as error function.

So loss function is calculated many times in a single training cycle but the cost function is calculated once.

cost = Σ loss

---

Logit Function:

The odds specify the ratio of prob. of success to that of failure.

$$\log\left(\frac{P(x)}{1 - P(x)}\right) = \beta_0 + \beta_1 . x$$

↓

logit ∿ odds
or log-odd

∴ in logistic regression, linear combination of inputs are mapped to the log (odds).

$$\text{sig}(x) = \frac{1}{1 + e^{-x}}$$

$$= \frac{1}{1 + \frac{1}{e^x}}$$

$$\boxed{\text{sig}(x) = \frac{e^x}{1 + e^x}}$$

$$P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

→ sigmoid function → gives a S-shaped curve. It always give a value of prob. ranging from (0,1).

# Estimation of parameters:

Unlike linear regression that use OLS, we use MLE.

There can be infinite sets of regression coeff. The Max. Likelihood Estimate is that set of regression coeff. for which the output prob. is max.

$$L(\beta, y) = \prod_{i=1}^{N} \left( \frac{P_i}{(1-P_i)} \right)^{y_i} \times (1-P_i)$$

(with "success" labeled over $\frac{P_i}{(1-P_i)}$)

For simplicity we use log likelihood.

---

## Performance:

### Confusion matrix

|  | | Predicted | |
|---|---|---|---|
|  |  | P | N |
| **actual** | P | TP | FN |
|  | N | FP | TN |

### Accuracy:

$$\frac{TP + TN}{TP + TN + FP + FN}$$

---

### Linear regression

$$y = mx + c$$

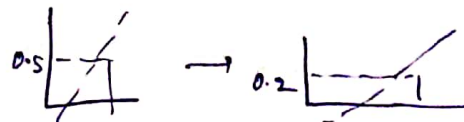where $m$ is slope, $c$ is intercept, and $x$ is input data points.

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$= \beta_0 + \beta_1 x$$

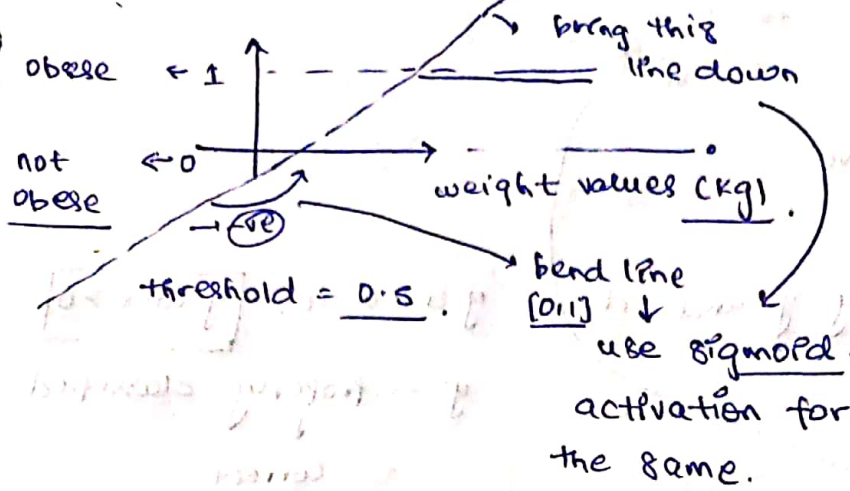$$y = \omega^T x + \theta$$

Linear regression is susceptible to outliers.

This will lead to changing thresholds. (0.5 → 0.2)

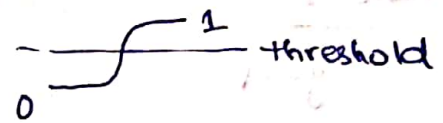Hence, we can't make use of MSE or OLS and also linear regression.

$\Rightarrow$ **conclusion**.

• ] $\rightarrow$ will always be classified as obese.

obese $\leftarrow 1$

not $\leftarrow 0$
obese

threshold = 0.5

$\rightarrow$ bring this line down

weight values (kg)

$\rightarrow$ bend line [0,1] $\downarrow$
use sigmoid activation for the same.

$\downarrow$
line keeps growing.

$\downarrow$
This fails miterably.

$\rightarrow$ squashed the straight line.

threshold

also used in deep learning too in the output dense layer.

$\text{sigmoid}(r) = \dfrac{1}{1+e^r}$
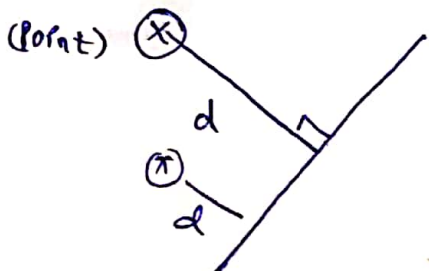
---

**assumptions:**

all points $\rightarrow$ +ve class or $c_1 \rightarrow +1$
" " $\rightarrow$ -ve " or $c_2 \rightarrow -1$ } initialization

$y = w^T x + b = 0$
$\rightarrow$ (passing through origin)

$\boxed{y = w^T \cdot x}$

points are linearly separable.

---

(point) $\otimes$

d

$\widehat{(\pi)}$
d

plane (w values) (coeff. values)

$d = \dfrac{w^T x + b}{\|w\|}$

if $w = $ unit vector
$\|w\| = 1$

line passing through origin $\Rightarrow b = 0$

sum of distances.

$\sum\limits_{i=1}^{n} w_i^T x_i$

$\boxed{d = w^T x}$

all distances will be +ve or zero or -ve depending on the side of the plane the point lies.

① case-1:

$$y = +1$$

above the plane distance = +ve

$$d = w^T x > 0$$

multiply $y$ and $d$

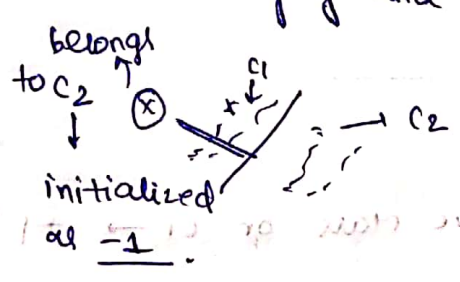$$y \cdot d > 0 \text{ ie. } \boxed{y \cdot w^T x > 0}$$

$y_i \to$ properly classified.
↓
correct.

② case-2:

$$y = -1$$

$$d = w^T x < 0$$

multiply $y$ and $d$.

-ve × -ve = +ve value

$$\boxed{y \cdot w^T \cdot x > 0}$$
↓
correctly classified.

③ case-3:

belongs to $c_2$
↓
initialized as -1

$$y = -1 \quad d = w^T x > 0 \text{ (above the plane)}.$$

multiply $\boxed{y \cdot w^T x < 0} \to$ incorrect classification.

④ case-4:

Ⓧ belongs to $c_1$

$$y = +1$$
$$d = -ve$$

$$\boxed{y \cdot d < 0} \to \text{ incorrect classification}.$$

cost function: $\max\left\{\sum\limits_{i=1}^{N} y_i \cdot w_i^T \cdot x_i\right\}$ → +ve → correct classif.

or

optimizer. → -ve → incorrect classif.

↓ given  ↓ given

Note: An outlier can change the entire game

update these coeff. unless I get max. value of cost function. ← compute these coefficients.

Solution:

let $z = w^T x$

$z = y \cdot w^T x$

$\max\{ \text{sigmoid}( y_i \underset{\uparrow}{\times} w_i^T \underset{\uparrow}{\times} x_i )\}$

operator

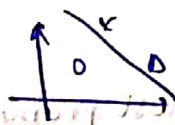sigmoid (z) = $\dfrac{1}{1+e^{-z}}$

after value calculation → value is transformed b/w 0 to +1.

Note:
Ill-Effect of an outlier is nullified.

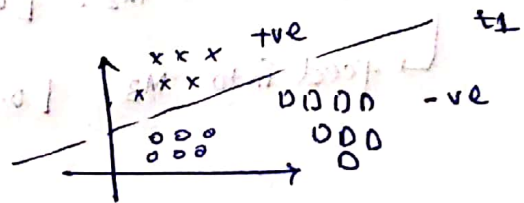Logistic Regression ( one vs Rest) or ( one vs All)

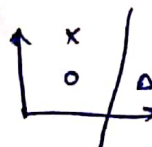Solving multiclass classification.
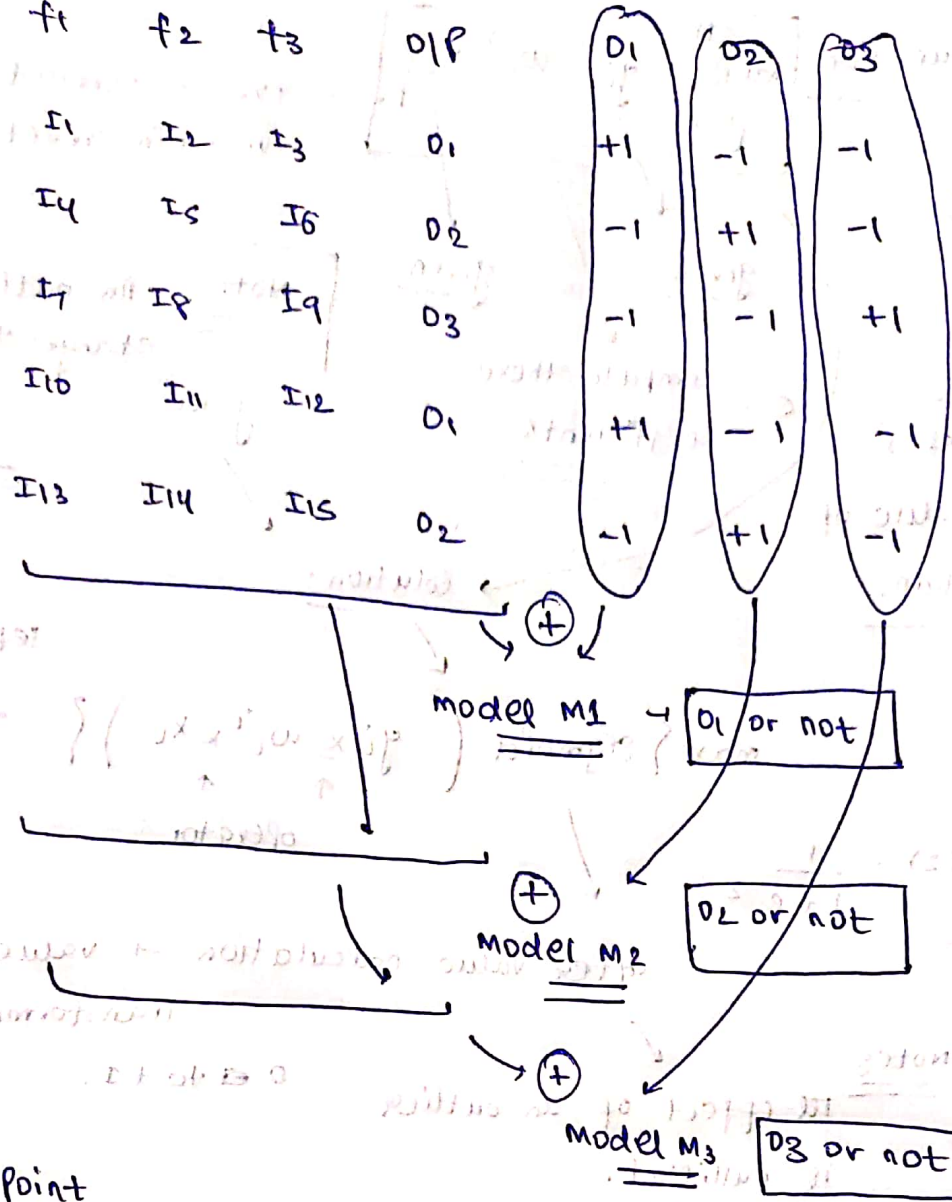
iteration 2
M₂ model

iteration 1 — M₁ model

x x x +ve
x x x
◻◻◻◻ -ve
○ ○ ○   ○○○
○ ○ ○   ○ ○ ○

+1

iteration 3  → M₃ Model

take one out and group rest all as 1.

Note:
take combinations.

Scanned by CamScanner

| | ft | f2 | t3 | O/P | O1 | O2 | O3 |
|---|----|----|----|-----|----|----|----|
| | $I_1$ | $I_2$ | $I_3$ | $O_1$ | +1 | -1 | -1 |
| | $I_4$ | $I_5$ | $I_6$ | $O_2$ | -1 | +1 | -1 |
| | $I_7$ | $I_8$ | $I_9$ | $O_3$ | -1 | -1 | +1 |
| | $I_{10}$ | $I_{11}$ | $I_{12}$ | $O_1$ | +1 | -1 | -1 |
| | $I_{13}$ | $I_{14}$ | $I_{15}$ | $O_2$ | -1 | +1 | -1 |

$\oplus$

model M1 → | $O_1$ or not |

$\oplus$

Model M2       | $O_2$ or not |

$\oplus$

Model M3   | $O_3$ or not |

New Point

     output Prob.

⤷ feed into M1: [0.20]

   ⤷ feed into M2: [0.25]      summation = 1

    ⤷ feed into M3: [0.55]

                      highest probability

overall → [0.20, 0.25, 0.55)      value

max.val = 0.55 → ans = | O3 |