## Logistic Regression

Input
→ Signals
→
→
→
→ Dendrites

Axon

nucleus

Output
Signals

net input

$z = \omega_0 + x_1\omega_1 + x_2\omega_2$

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

### Sigmoid Function

1·0

0·5

0·0

−8  −6  −4  −2   0   2   4   6   8

→ ON

OFF

$1$   $\omega_0$

$x_1$   $\omega_1$

$Z$   $\sigma(z)$

$\omega_2$

$x_2$

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

in reality biological neurons
dont just have 0 or 1
They have continuous
like
<u>Sigmoid function</u>
it imitates the on/off values

Activation.

$\varphi(z) > 0.5 \rightarrow 1$

$\varphi(z) < 0.5 \rightarrow 0$

→ Logistic Neuron can be used as
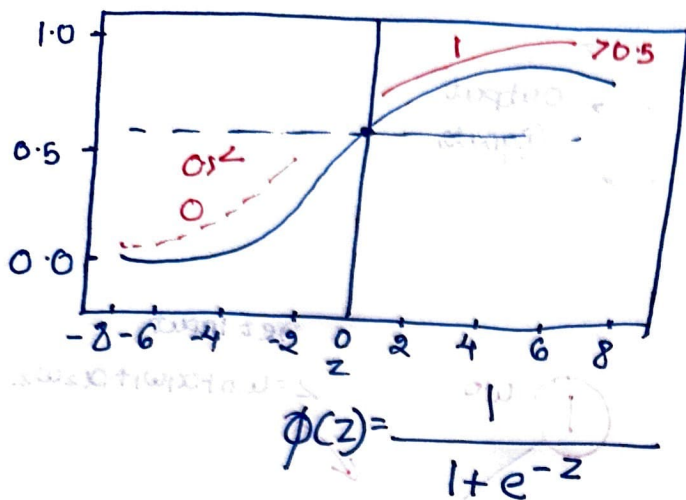a binary classifier

Probability → The input has a
Property

The input is a dog/not
human/dog
Input has a property

## Properties of Logistic Neuron :-



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

What points are on the fence?
$$\phi(z) = 0.5, \quad z = 0$$
on the boundary.

$$\omega_0 + \omega_1 x_1 + \omega_2 x_2 = 0$$
$$\longrightarrow \omega_0 + x \omega^T = 0$$
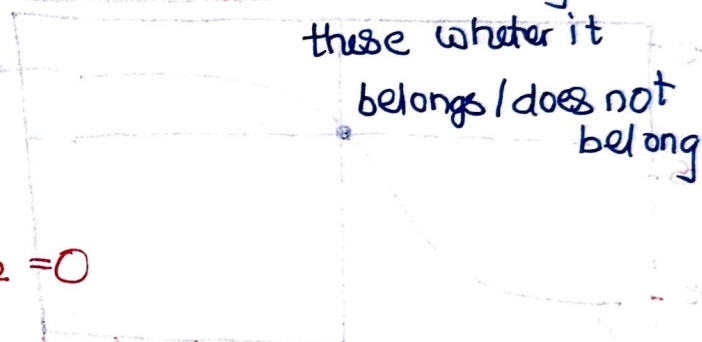$$\text{(in higher dimensions)} \Big] \text{hyperplane}$$

## Binary Classifier / Linear Classifier

Logistic Neuron
(Sigmoid fn)
can be used as a binary classifier.

$$\phi(z) \xrightarrow{\quad} 0.5 \rightarrow 1 \quad \text{belong to class}$$
$$\xrightarrow{\quad} 0.5 \rightarrow 0 \quad \text{not belongs to class}$$

what is the decision boundary b/w these wheter it belongs / does not belong

② why we look into the output of logistic neurno / sigmoid function and interpret it as probability?

$$P \longrightarrow \text{odds}(P) = \frac{P \quad \leftarrow \text{probability}}{1-P \quad \leftarrow \text{(complement)}} \rightarrow \log \text{it} \, (P)$$
$$[0,1] \qquad [0,\infty] \qquad \qquad = \log\left(\frac{P}{1-P}\right)$$
$$[-\infty, \infty]$$

$$Z : \text{logit of some } P$$
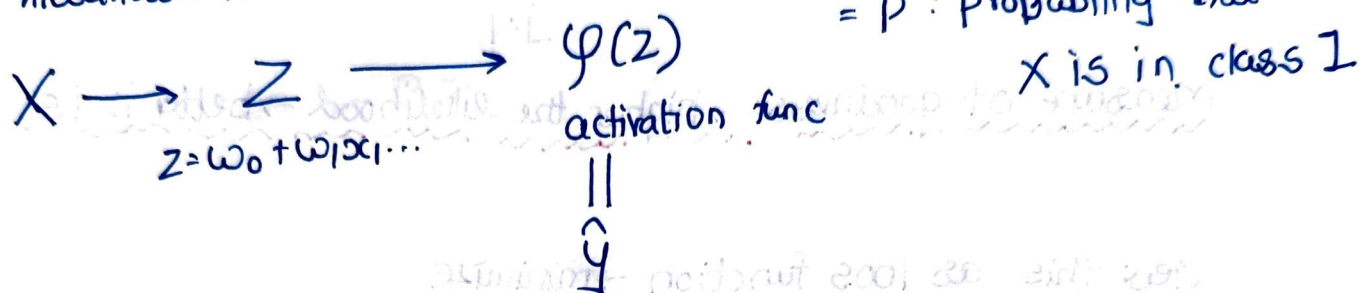$$\log \frac{P}{1-P} = Z \implies P = \frac{1}{1 + \exp(-z)}$$

# The loss function (negative log likelihood)

Data
↑
$$L(\underset{\underset{\text{label}}{\downarrow}}{x}, \underset{\underset{\text{weights}}{\searrow}}{y}, w)$$

→ Find $w$ that minimizes the Loss function.

---

likelihood $x$

$$X \longrightarrow Z \longrightarrow \underset{\text{activation func}}{\varphi(Z)} \qquad = P : \text{probability that } X \text{ is in class 1}$$

$$z = w_0 + w_1 x_1 \cdots$$

$$\| \\ \hat{y}$$

$X →$ true label $y$

→ predicted → $\hat{y}$

probability $= 0.8$

80% good

but 100% would be better.

* $y →$ true label

* $\hat{y} → P →$ predicted (probability)

$y_r = 1, P = 0.8$

$y = 0, P = 0.2$ I am 80% good as per the complement

$$\text{likelihood } (x) = \hat{y}^{y} (1-\hat{y})^{(1-y)}$$

when $y = 1$ ~~const~~ $\hat{y} = 0.8$

$$\text{likelihood } (x) = \hat{y}^{1} (1-\hat{y})^{(1-1)} \qquad \hat{y} = 0.8$$

$$= 0.8^{1}$$

$$y = 0, \hat{y} = 0.2$$

$$\text{likeli} = \hat{y}^{0} (1-\hat{y})^{(1-0)}$$

$$= 0.2^{0} (1-0.2)^{(1-0)} = 0.8$$

$$X = \{ x_1, \ldots\ldots, x_n \}$$

$$\text{likelihood } (X) = \text{product of individul likelihood}$$

$$= \prod_{j=1}^{n} \text{likelihood } (x_j)$$

measure of goodness → higher the likelihood → better it is.

View this as loos function → minimize
likelihood → maximize

$$Lh(X) \xrightarrow{\text{max}} Log(Lh(x)) = \sum_{j=1}^{n} log(\text{likelihood } (x_j))$$

maximize the log (Lh) insted of Lh directly.

Convert to
{ minimize

* Loss function :- nll (X) = $-\sum_{j=1}^{n} log(\text{likelihood } (x_j))$

So if we minimize loss we maximize likelihood.

# Logistic Regression and Regularization

labels $\to \{0,1\}$

function $\omega =$ Logistic Regression - Stochastic Fit $(X, y)$

Initialize randomly a d-dimensional vector $\omega$ and a scalar b

Shuffle the rows (points) of $X \to$ Shuffle

for $j=1$ to n epochs

   for $t=1$ to $n/k$ $\longrightarrow$ hyperparameter $k$ (size of mini batch)

                             one update per batch.

$$b = b - \beta \sum_{j=(t-1)k+1}^{tk} \left( \underset{\hat{y}_i}{\sigma(b+X_j\omega^T)} - \underset{-y_i}{y_i} \right)$$

        $\leftarrow$ sigmoid fn

    for $k=1$ to d:

$$\omega_k = \omega_k - \beta \sum_{j=(t-1)k+1}^{tk} \left( \sigma(b+X_j\omega^T) - y_j \right) X_{j,k}$$
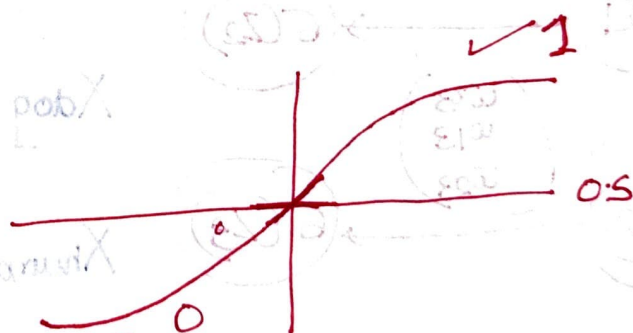
return $\omega, b$

$L(X, y, \omega):$ $\omega$ that minimizes loss

gradient descent

when $k=1$: identical to perceptron

difference is $\sigma \to$ sigmoid
is used for
calculation

SGD $\to$ Logistic neuron

$X\omega^T + b$

$\to z$

$\to \sigma$

*From one vs rest (from binary to multilabel classification):

$\mathcal{L}(X, y, \omega):$ ~~$\omega$ that minimizes the loss function~~



$$z = \omega_0 + \omega_1 x_1 + \omega_2 x_2$$

logistic neuron

$$\sigma(z) = \frac{1}{1+e^{-z}} \Big\} \text{ Prob } X \text{ is in class 1}$$

CAT / DOG / HUMANS : $X = [X_{CAT}, X_{DOG}, X_{CAT}]$

$$\omega^{(1)} = \begin{pmatrix} \omega_{01} \\ \omega_{11} \\ \omega_{21} \end{pmatrix}$$

$$\omega^{(2)} = \begin{pmatrix} \omega_{02} \\ \omega_{12} \\ \omega_{22} \end{pmatrix}$$

$$\begin{pmatrix} \omega_{03} \\ \omega_{13} \\ \omega_{23} \end{pmatrix}$$

Train 3 classifiers

$X_{cat}:$  $\big[ X_{cat}, \{X_{dog}, X_{human}\} \big]$  0.3 cat
              1          0

$X_{dog}:$  $\big[ X_{dog}, \{X_{human}, X_{cat}\} \big]$  0.6 dog
              1          0

$X_{human}$  $\big[ X_{human}, \{X_{cat}, X_{dog}\} \big]$  0.8 human
              1          0

$$\frac{0.3}{1.7} + \frac{0.6}{1.7} + \frac{0.8}{1.7}$$

→ output prob distribution

Final Classification = $\max\{0.3, 0.6, 0.8\}$

0.8 → human (highest probability)

$\big(0.3 + 0.6 + 0.8$ don't sum up to 1$\big)$ but its fine $\Big\} \dfrac{\sigma(z_i)}{\sum_i \sigma(z_i)}$
             = 1.7

if we have k labels :   k classifiers (X)

every classifier uses entire
dataset but with
different labels.

~~fooooooooo~~ $\left(\dfrac{k(k-1)}{2}\right)$   $\{X_{CAT}, X_{dog}\}$