

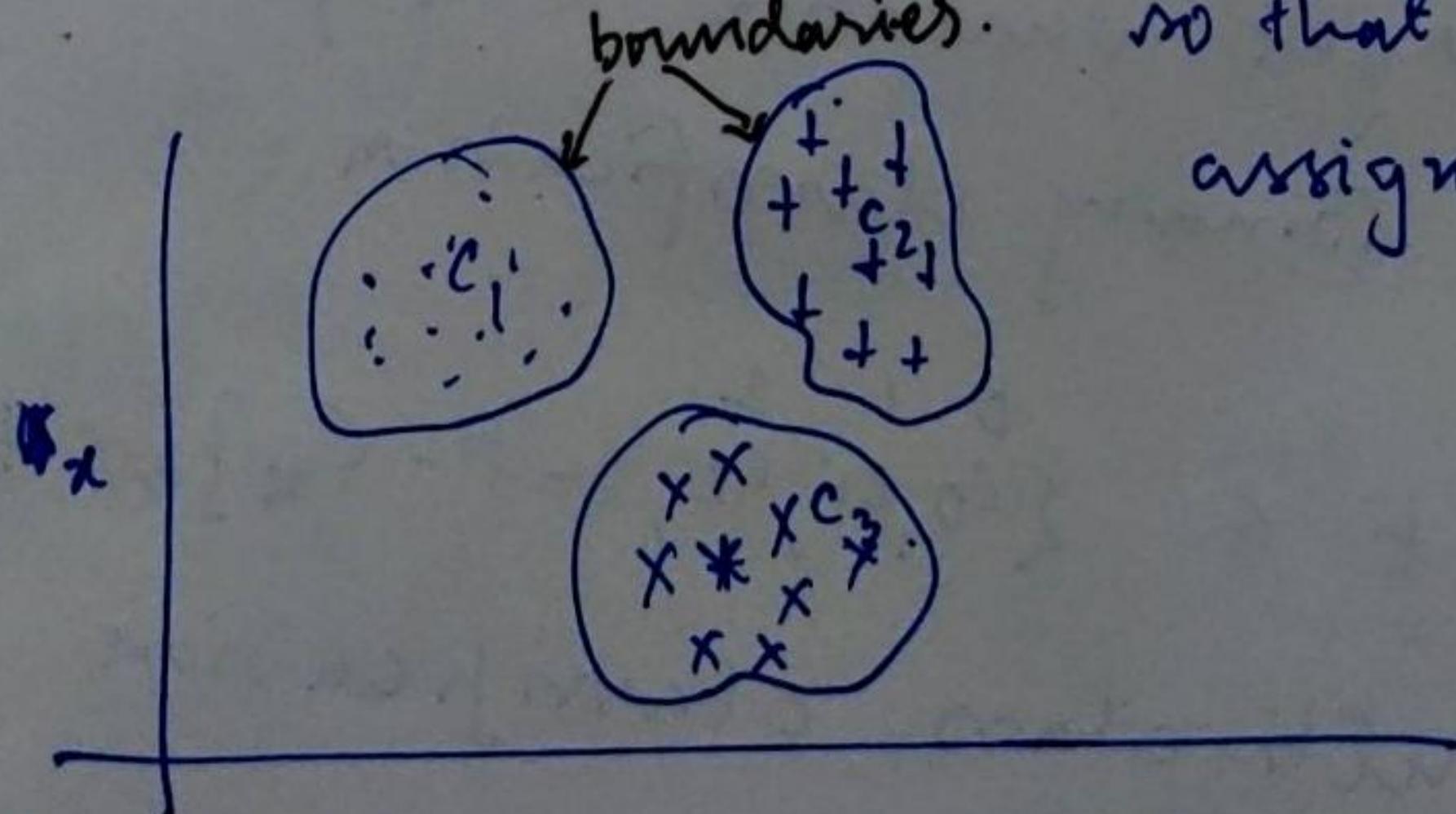
"Classification and logistic regression"

— Classification prob are just like regression prob, except that the values of target variable y , which we want to predict takes only a small number of discrete values. ($\text{Ans } k = \{1, 2, \dots, K\}$)

→ The goal in classification prob is to take an input vector \underline{x} and to assign it to one K discrete classes C_k , where

$$K = \{1, 2, \dots, K\}$$

→ Common sener scenario classes are disjoint boundaries. so that each input is assigned only one class.



— Input space is thereby divided into Decision surface/ decision boundary.

Input space is divided into decision spaces, whose boundaries are decision boundaries.

- In general, is $\underline{x} \in \mathbb{R}^{D \times n}$ (D dimensional,

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix} \in \mathbb{R}^{D \times n} \quad n: \dim(\text{features})$$

If linear models of classification is adopted, which means "decision surfaces" are linear functions of the input vector \underline{x} and hence are defined by $(\mathbb{R}^D)^{(n-1)}$ -dimensional hyperplane.

Binary & Multiclass classification.

The target variable t or $y \in \{0, 1\}$ is called Binary classification

If $t \in \{c_0, c_1, c_2, \dots, c_k\}$ then it is multi-class classification.

$$t = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix},$$

$$t = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix},$$

$$t = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Three distinct approaches for classification prob.

Deterministic: This involves constructing a "discriminant function", which assigns each input vector \underline{x} a class.

probabilistic approach:

- Is to model the conditional probability $p(C_k | \underline{x})$

- two different approaches to model $p(C_k | \underline{x})$.

A) Model them directly. Represent them as parametric models and then optimizing the parameters using a training set.

B) Second, is generative approach. in which we model the class-conditional densities, given by $p(\underline{x} | C_k)$, together with prior probabilities $p(C_k)$ for the classes. and.

using Baye's theorem. = $P(\underline{x}, c_k) / P(\underline{x})$

$$P(c_k | \underline{x}) = \frac{P(\underline{x} | c_k) P(c_k)}{P(\underline{x})}$$

c_k are the classes

$$P(\underline{x}) = \sum_{k=1}^K P(\underline{x} | c_k) P(c_k)$$

We model the joint distribution in generative model.

Classification: Binary ($K=2$). $\stackrel{c_0}{\sim} \stackrel{c_1}{\sim}$

e.g., Email : Spam / Not Spam

Online transaction: Fraudulent (Yes/No)

Tumor : Malignant / Benign

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \longrightarrow y^{\text{act}} \in \{0, 1\}$$

)
0: negative class
1: positive class.

Feature vector

example: let's take one feature x_1 : "tumor size".

$y \in \{0, 1\}$
"Not malignant" \swarrow "Yes!" \searrow
/Benign Malignant

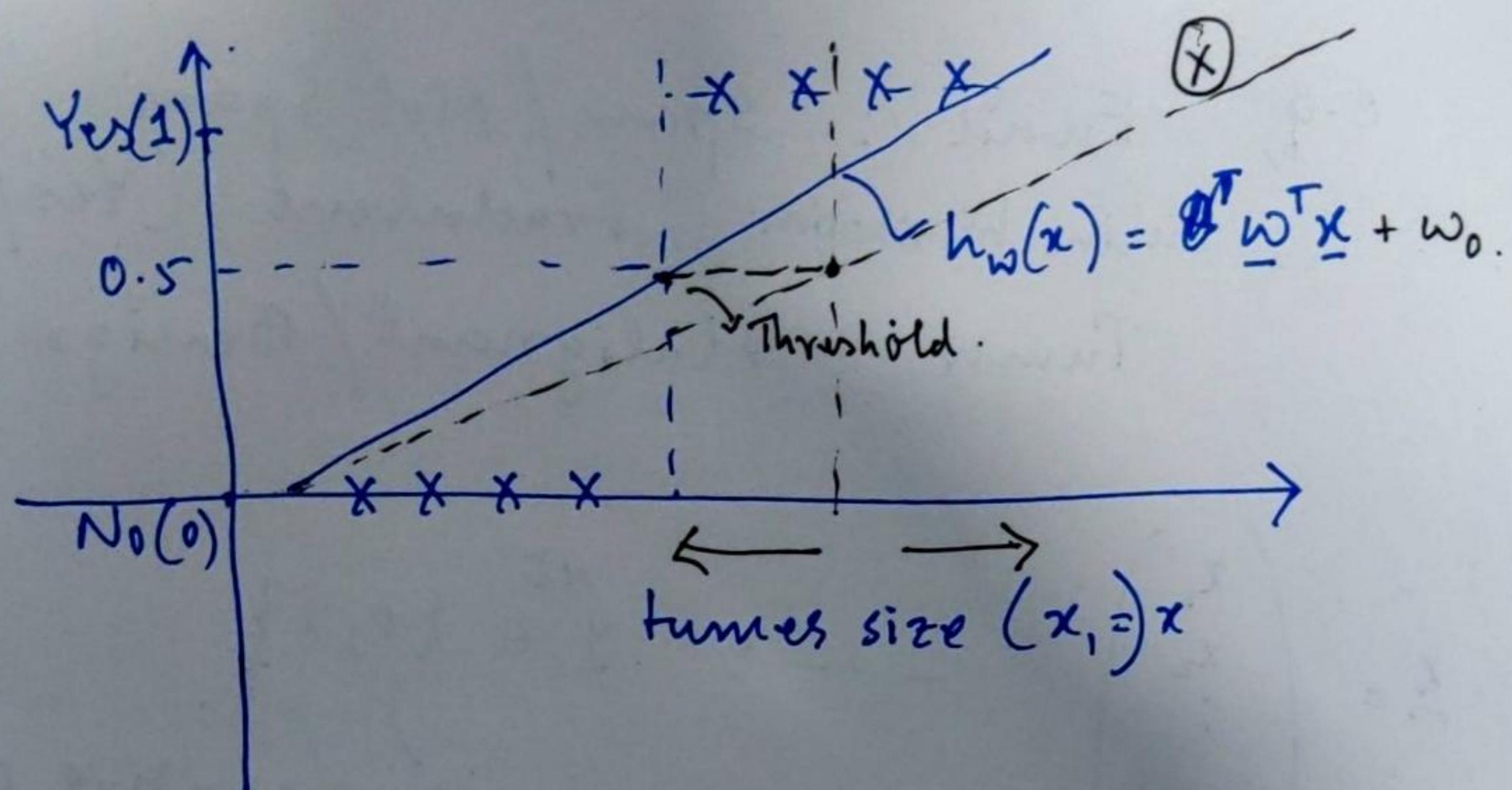
usual approach may be:

1) Form the hypothesis function, and.

$$h_w(\underline{x}) = \underline{w}^T \underline{x} + w_0 = \underline{w}^T \underline{x} + w_0$$

Find the parameters, assuming $h_w(\underline{x})$ is a continuous function.

how would we classify then.



define a threshold classifier output
 $h_w(x) = 0.5$, ($h_w(x)$ is found by fitting the data)

Data

$(x_1^{(1)}, 1)$	\leftrightarrow	$(x_1^{(1)}, y)$	if $h_w(x) \geq 0.5$, assign $y = 1$
$(x_1^{(2)}, 0)$	\leftarrow	$(x_1^{(2)}, N)$	if $h_w(x) < 0.5$, assign $y = 0$.
$(x_1^{(3)}, 0)$	\leftarrow	$(x_1^{(3)}, N)$	
\vdots			
$(x_1^{(n)}, 1)$	\leftarrow	$(x_1^{(n)}, y)$	please Note x is a continuous variable.

2 we are assuming $h_w(x) = w^T x + w_0$ is Conti

P Issues: Although final target variable
 $y_{\text{wt}} = 0 \text{ or } 1$

but, $h_w(x)$ can be > 1 or < 0 .

which DOES NOT make sense (Intuitively)

This introduce, nonlinear boundedness:

As y is a discrete variable (for class identification) we change the hypothesis function definition by introducing a non-linear functional map.

logistic regression $\rightarrow h_w(x) = g(\underline{w}^T \underline{x})$

$$= g(\underline{w}^T \underline{x} + w_0)$$

so that $0 \leq g(\underline{w}^T \underline{x}) \leq 1$

$\therefore 0 \leq h_w(x) \leq 1$

$\underline{x}' = \begin{pmatrix} 1 \\ \underline{x} \end{pmatrix}$

$\underline{w} = \begin{pmatrix} w_0 \\ \underline{w} \end{pmatrix}$

Bounded hypothesis.

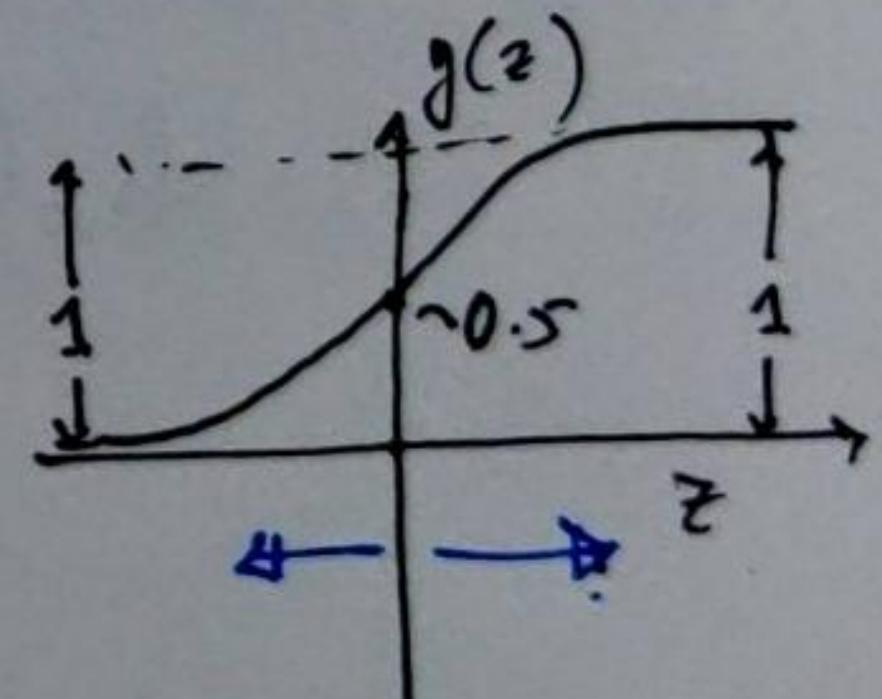
$g(\cdot)$: is called activation function.

g^{-1} : is called Link function. $\tilde{\underline{x}} = \begin{pmatrix} x_0 \\ \underline{x}_1 \end{pmatrix}$

$$\tilde{\underline{w}} = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix}$$

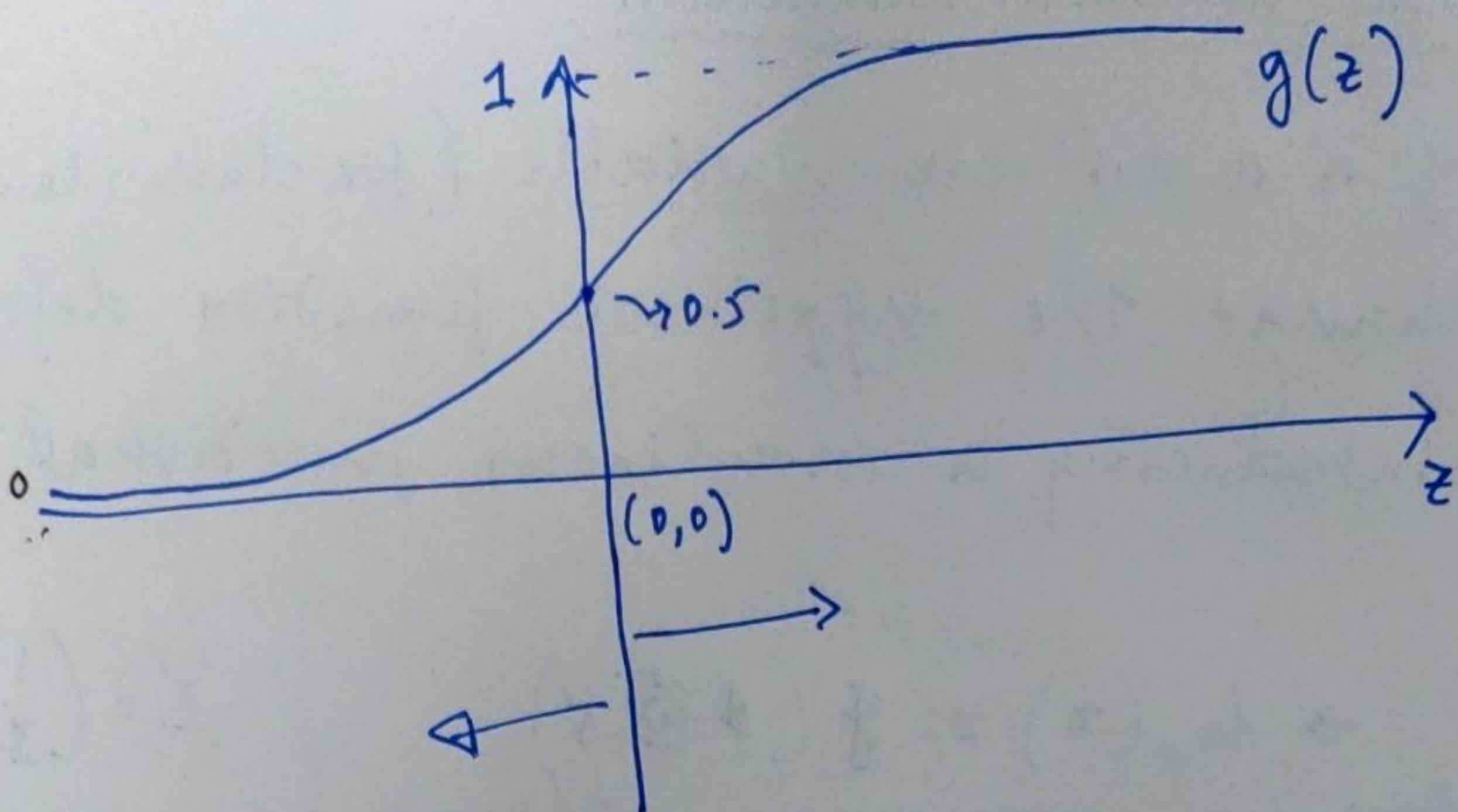
One type of activation function is "Sigmoid Function"

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



hence, for logistic regression, the hypothesis function

$$h_w(x) = g(\underline{w}^T \underline{x}) = \left(\frac{1}{1 + e^{-\underline{w}^T \underline{x}}} \right)$$



Notice that $g(z) \rightarrow 1$ as $z \rightarrow \infty$
 and $g(z) \rightarrow 0$ as $z \rightarrow -\infty$ hence.
 $h_w(x)$ is always bounded betⁿ 0 to 1

Interpretation of hypothesis outcome. / A probabilistic viewpoint.

$h_w(x)$ *: estimated probability that $y = 1$
 on input x

$$\equiv P(y=1 \mid x; w)$$

Example:

$$x = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ \text{tumor size} \end{pmatrix}$$

so if $h_w(x) = 0.7 \Rightarrow$ tell patient 70% chance of tumor being malignant.

$0 \leq h_w(\underline{x}) \leq 1 \Rightarrow$ can be assigned to a prob.

$$h_w(\underline{x}) = P(y=c_k \mid \underline{x}; \underline{w}).$$

↑ general. posterior probability.

probability that $y=1$, given \underline{x} , parameterized by \underline{w}

$$\Rightarrow P(y=0) + P(y=1) = 1.$$

or, $\Rightarrow P(y=0 \mid \underline{x}; \underline{w}) + P(y=1 \mid \underline{x}; \underline{w}) = 1 \quad \}^*$

Decision boundary

The goal in classification is to take an input vector \underline{x} and to assign it to one of K discrete classes c_k where

$$k=1, \dots, K$$

In most common scenario, the classes are taken to be disjoint, so that each input is assigned to only one and only one class.

Input space is thereby divided into decision regions whose boundaries are called decision boundaries

or decision surfaces.

Linear models of classification, involves decision surfaces. are linear functions of input variable \underline{x}

$$\begin{array}{c}
 h_w(\underline{x}) \\
 \downarrow \\
 \text{if } h_w(\underline{x}) = 0, \text{ then } \{\underline{x}\} \\
 \text{if } h_w(\underline{x}) = 1, \text{ then } \{\underline{x}\} \\
 \downarrow \\
 \underline{x} \in \mathbb{R}^D \quad \underline{w}^T \underline{x} \in \mathbb{R} \\
 \downarrow \\
 \alpha(\underline{x})
 \end{array}
 \quad
 \begin{aligned}
 h_w(\underline{x}) : \mathbb{R}^D &\longrightarrow \{0, 1\} \\
 \Rightarrow h_w(\underline{x}) : \mathbb{R}^D &\longrightarrow \{C_1, C_2\} \\
 \Rightarrow g(\underbrace{w_0 + w_1 x_1 + \dots + w_D x_D}_{\text{linear function}}) : \mathbb{R}^D &\longrightarrow \{C_1, C_2\}.
 \end{aligned}$$

$w_0 + w_1 x_1 + \dots + w_D x_D$ is a $D-1$ dimensional hyperplane. $\in \mathbb{R}^D$.

The decision surfaces are defined as.

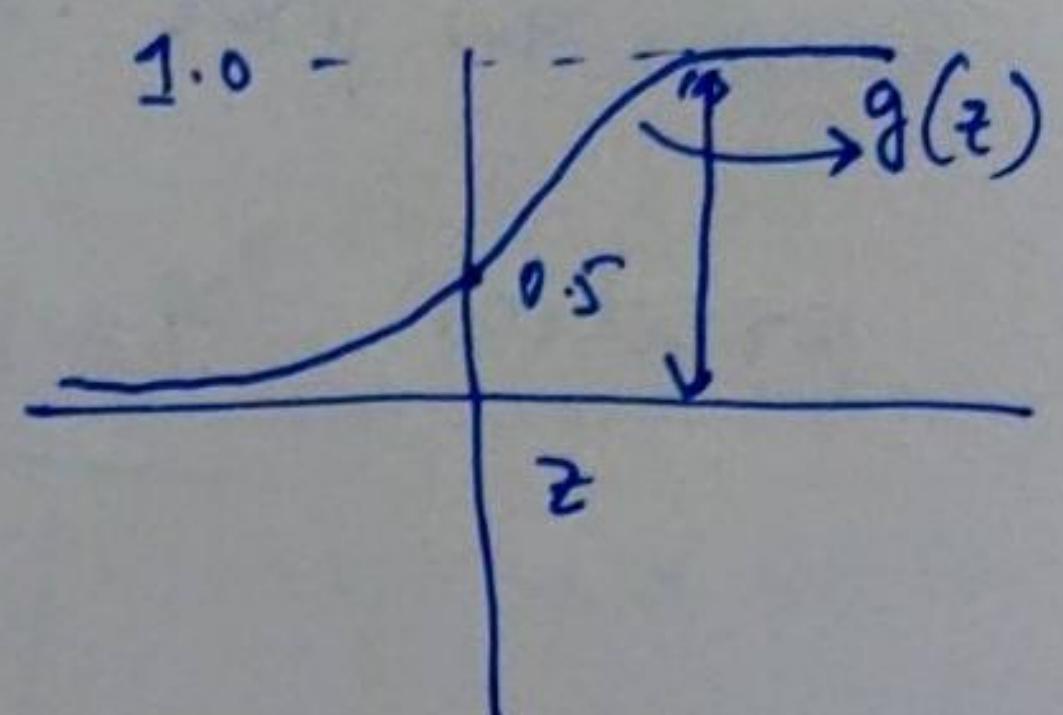
$$\begin{array}{c}
 g(w_0 + w_1 x_1 + \dots + w_D x_D) = \text{const.} \\
 \text{decision surface} \rightarrow w_0 + w_1 x_1 + \dots + w_D x_D = \text{const}
 \end{array}$$

hence, decision surfaces are linear functions of \underline{x} , even if the function g is non-linear.

$h_w(\underline{x}) = g(\underline{w}^T \underline{x})$: called generalized linear model

Example 3: linear case: Decision boundary

$$h_w(\underline{x}) = g(\underline{w}^T \underline{x}) \\ = g(w_0 + w_1 x_1 + w_2 x_2)$$



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

If $\xrightarrow{\text{Case I}} y=1$ for $h_w(\underline{x}) > 0.5$ # threshold definition

$$\Rightarrow g(\underline{w}^T \underline{x}) > 0.5$$

$$\Rightarrow \underline{w}^T \underline{x} > 0. \quad [\text{From the graph}]$$

$\xrightarrow{\text{Case II}}$ decision surface/boundary

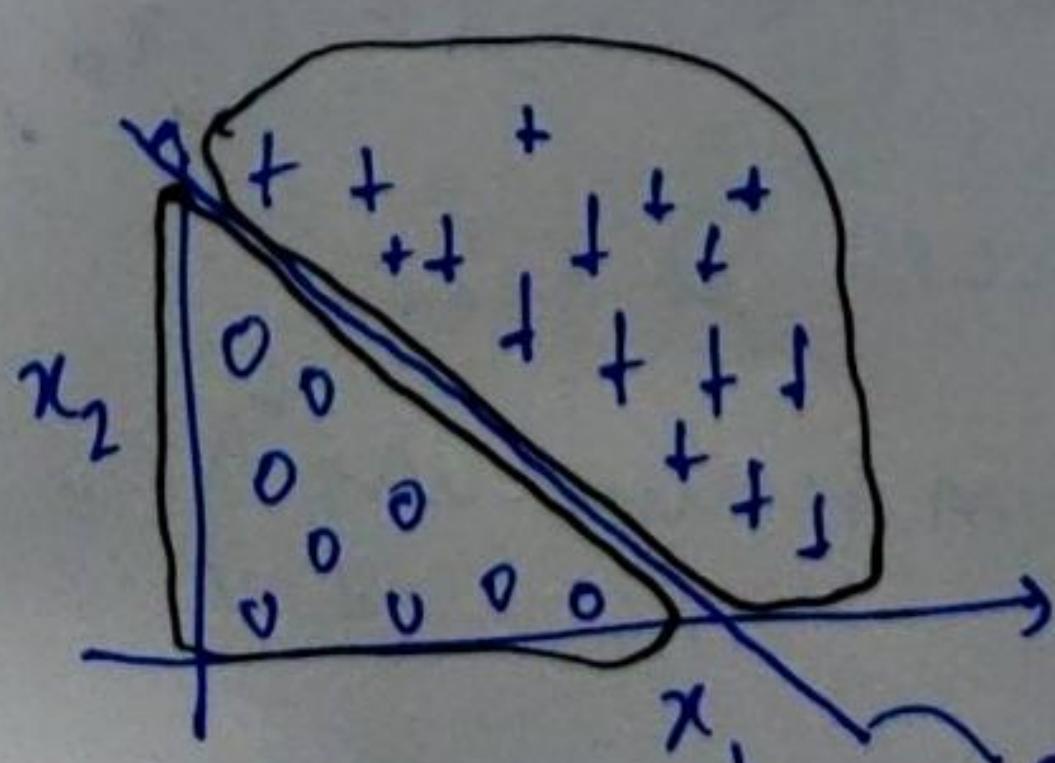
Case II If $y=0$ for $h_w(\underline{x}) \leq 0.5$

$$\Rightarrow g(\underline{w}^T \underline{x}) < 0.5$$

$$\Rightarrow z = \underline{w}^T \underline{x} < 0. \quad [\text{From graph}]$$

decision surface/boundary in 2D

Example 4



$$\underline{w} = \begin{pmatrix} -3 \\ 1 \\ 1 \end{pmatrix} \text{ & let's assume weights.}$$

$$h_w(\underline{x}) = g(w_0 + w_1 x_1 + w_2 x_2)$$

decision boundary (line in 2D)

As per the linear model of classification.

predict " $y = 1$ " if $\underbrace{-3 + x_1 + x_2 \geq 0}_{\bullet w^T x}$

$x_1 + x_2 \geq 3$: Eq of decision boundary

$y = 0$ if $x_1 + x_2 < 3$.

Till now we talked about linear transformation of input variable \underline{x} .

If the classes can be separated by linear transformation of the input variables \underline{x} , then these classes are linearly separable. & decision boundaries are hyperplanes of $\text{dim} = D - 1$

previous example $\underline{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \in \mathbb{R}^2$.

hence $D = 2$.

$\text{dim}(\text{decision boundary}) = 2 - 1 = 1$
line ↗

Non-linear decision boundary

The "decision boundary" to separate classes can be done if non-linear transformations are done on the input variables, using a vector of basis functions $\phi(\underline{x})$.

$$h_w(\underline{x}) = g \left(w_0 + w_1 \underline{x}_1 + w_2 \underline{x}_2 + w_3 \underline{x}_1^2 + w_4 \underline{x}_2^2 \right)$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$

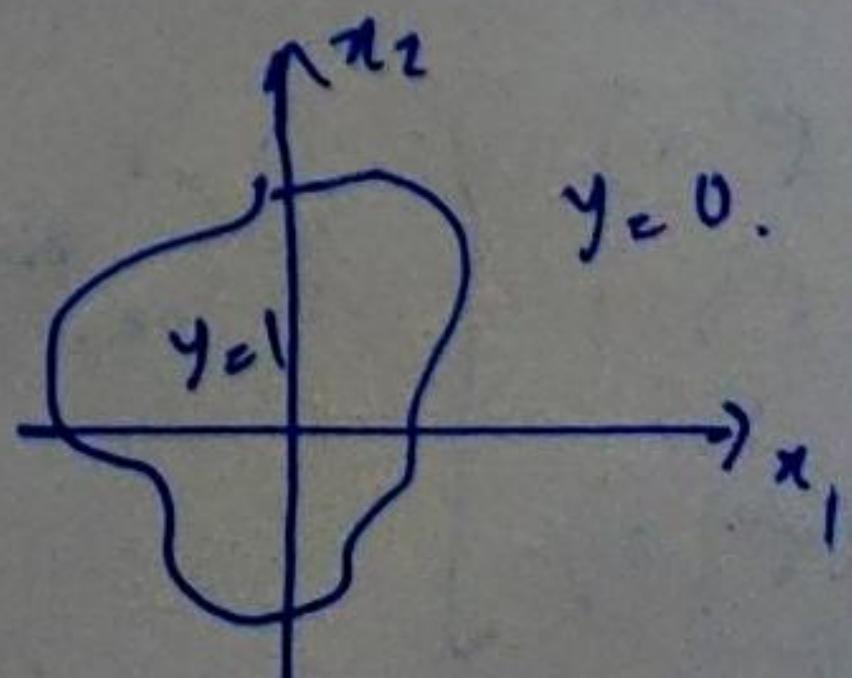
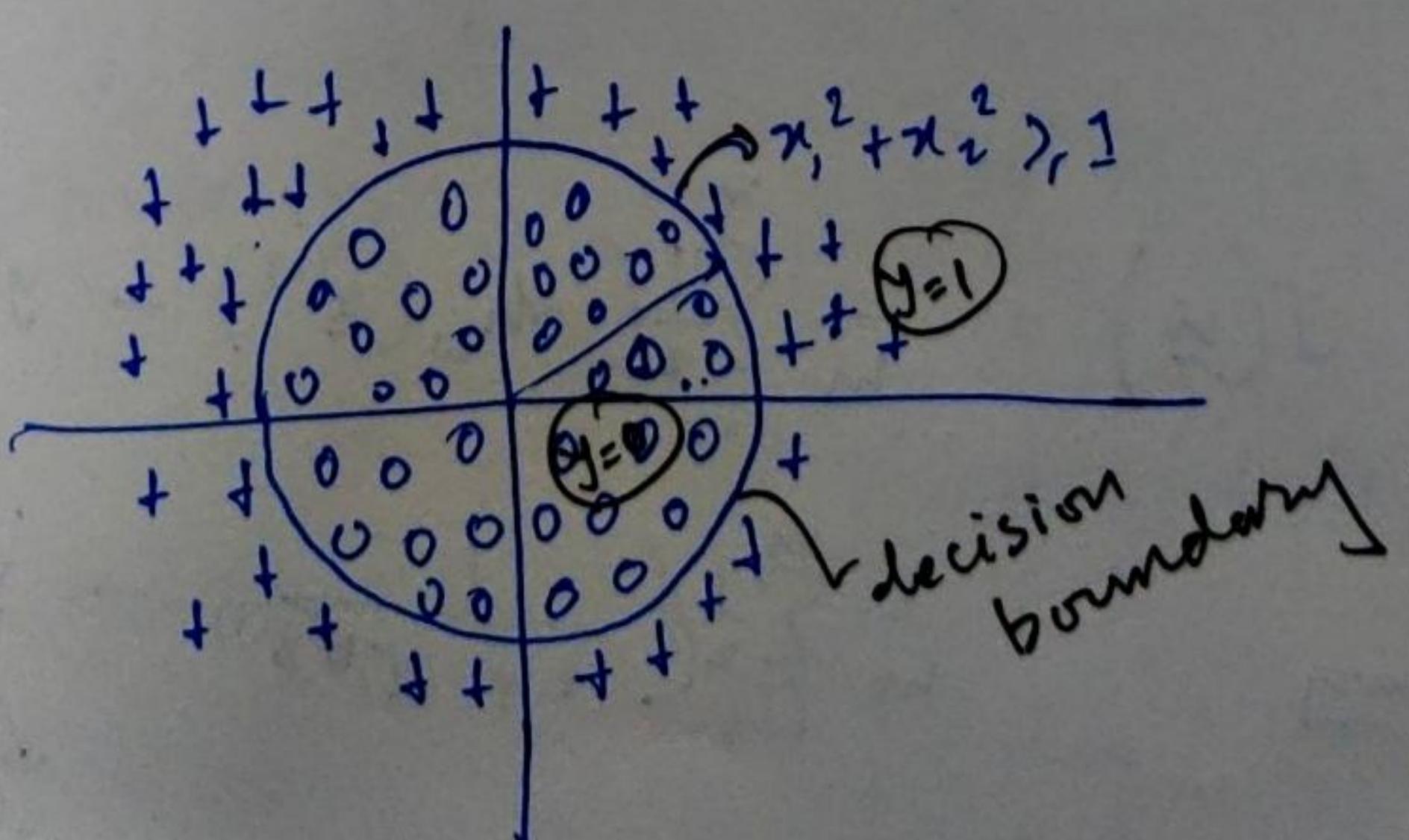
$\phi_1(\underline{x}_1) \quad \phi_2(\underline{x}_2) \quad \phi_3(\underline{x}_1) \quad \phi_4(\underline{x}_2)$

If weight vector is $w = \begin{pmatrix} -1 \\ 0 \\ 0 \\ 1 \end{pmatrix}$ then known weights.

$$h_w(\underline{x}) = g(-1 + 0 \cdot \underline{x}_1 + 0 \cdot \underline{x}_2 + 1 \cdot \underline{x}_1^2 + 1 \cdot \underline{x}_2^2)$$

predict "y=1" if $g(-1 + \underline{x}_1^2 + \underline{x}_2^2) > 0 \Rightarrow -1 + \underline{x}_1^2 + \underline{x}_2^2 > 0$

$$\Rightarrow \underline{x}_1^2 + \underline{x}_2^2 > 1$$



logistic regression : Cost function : How to find the weights??

Given the training set (m -data set)

$$\{(x^1, y^1), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$$

and,

$$\underline{x} = \begin{pmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ \vdots \\ x^{(m)} \end{pmatrix} \in \mathbb{R}^{n+1}$$

$$\underline{x} = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n+1}$$

$$x_0 = 1, \quad y \in \{0, 1\}$$

$$h_w(\underline{x}) = \left(\frac{1}{1 + e^{-w^T \underline{x}}} \right)$$

How to choose the parameters of the model?

Now cost function can not be defined as least square (like) for regression; for discriminant function $h_w(\underline{x})$ for ~~practical~~ practical reason.

$$J(w) = \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

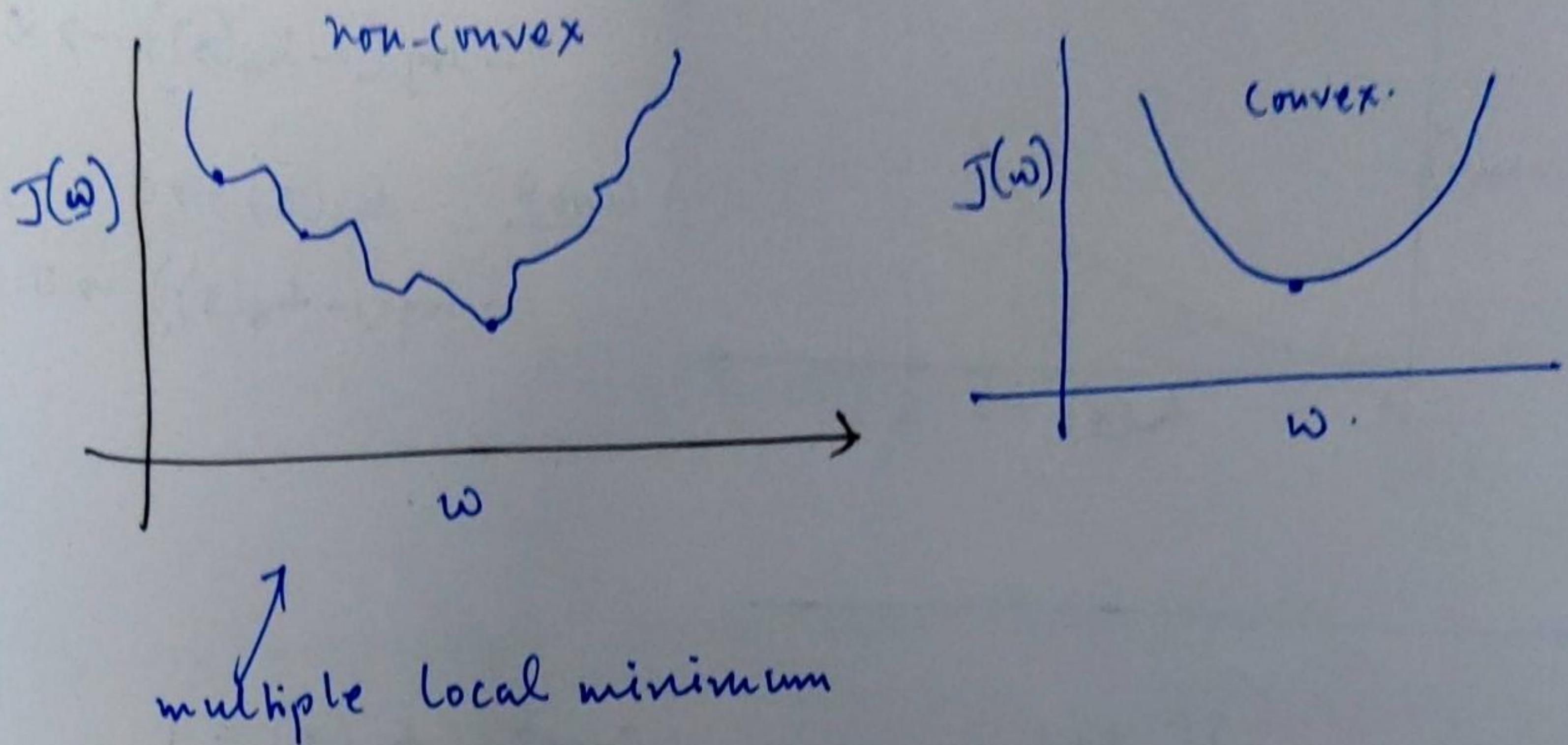
Data.

$$\left. \begin{array}{l} (x^{(1)}, 1) \\ (x^{(2)}, 1) \\ \vdots \\ (x^{(m)}, 0) \end{array} \right\} \quad \begin{array}{l} y \sim \\ 0, 1 \end{array} \text{ Binary} \quad = \frac{1}{m} \sum_{i=1}^m \left(\frac{1}{1 + e^{-w^T x^{(i)}}} - y^{(i)} \right)^2$$

$y^{(i)} \in \{0, 1\}$

generates a non-convex function.

→ Difficult to optimize (Non-convex)



(working at the individual term in the loss function $J(w)$, for a particular i)

$$\text{Cost}(h_w(x) - y) = \frac{1}{2} (h_w(x) - y)^2$$

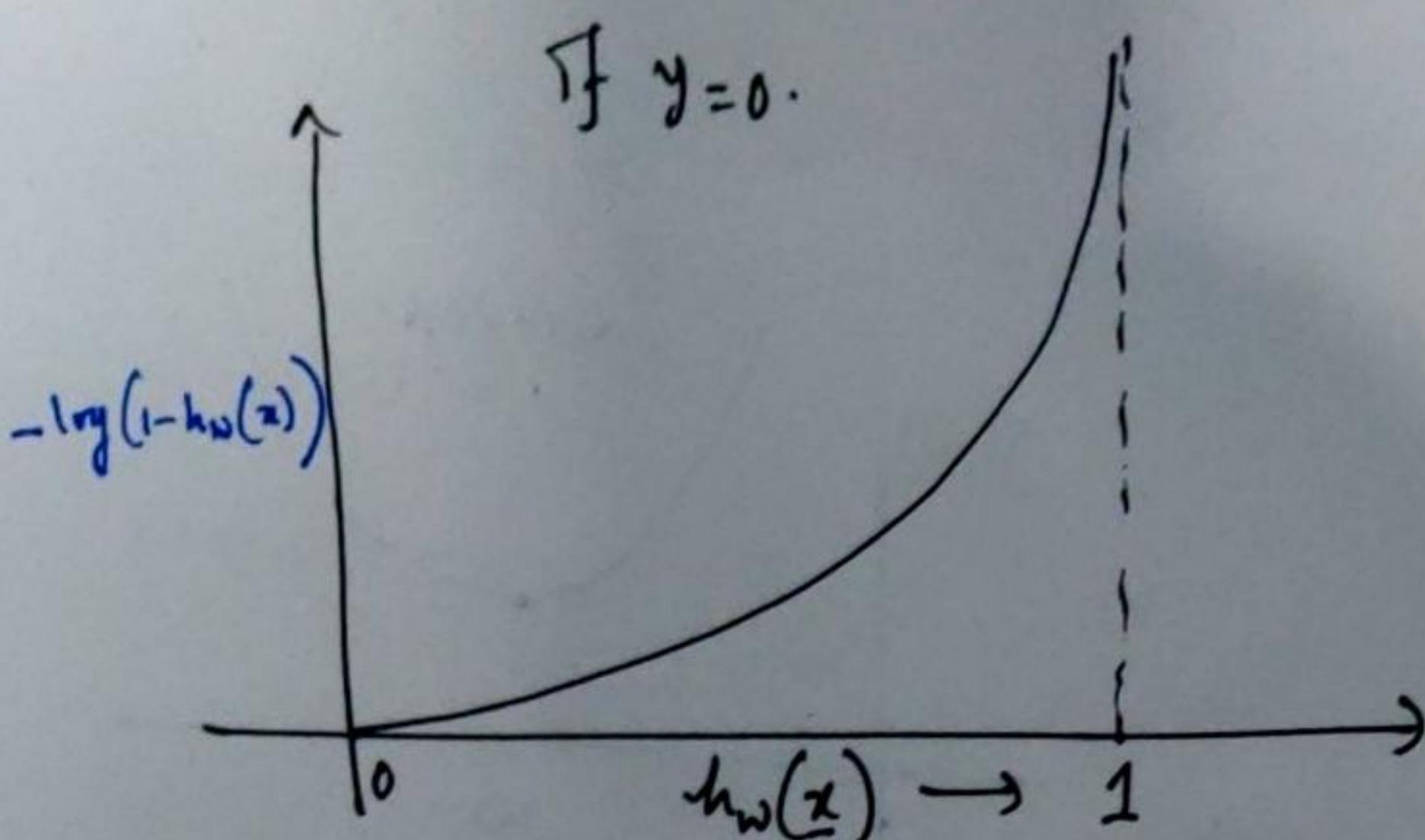
Not working.

Need to look for an alternative function.

→ An alternative is to convert it into a monotonic function. - "log" of the squared loss.

$$\text{Cost}(h_w(x) - y) = \begin{cases} -\log(h_w(x)) & \text{if } y=0 \\ -\log(h_w(x)-1) & \text{if } y=1 \\ 1-h_w(x) & \text{else} \end{cases}$$

Analyzing The cost



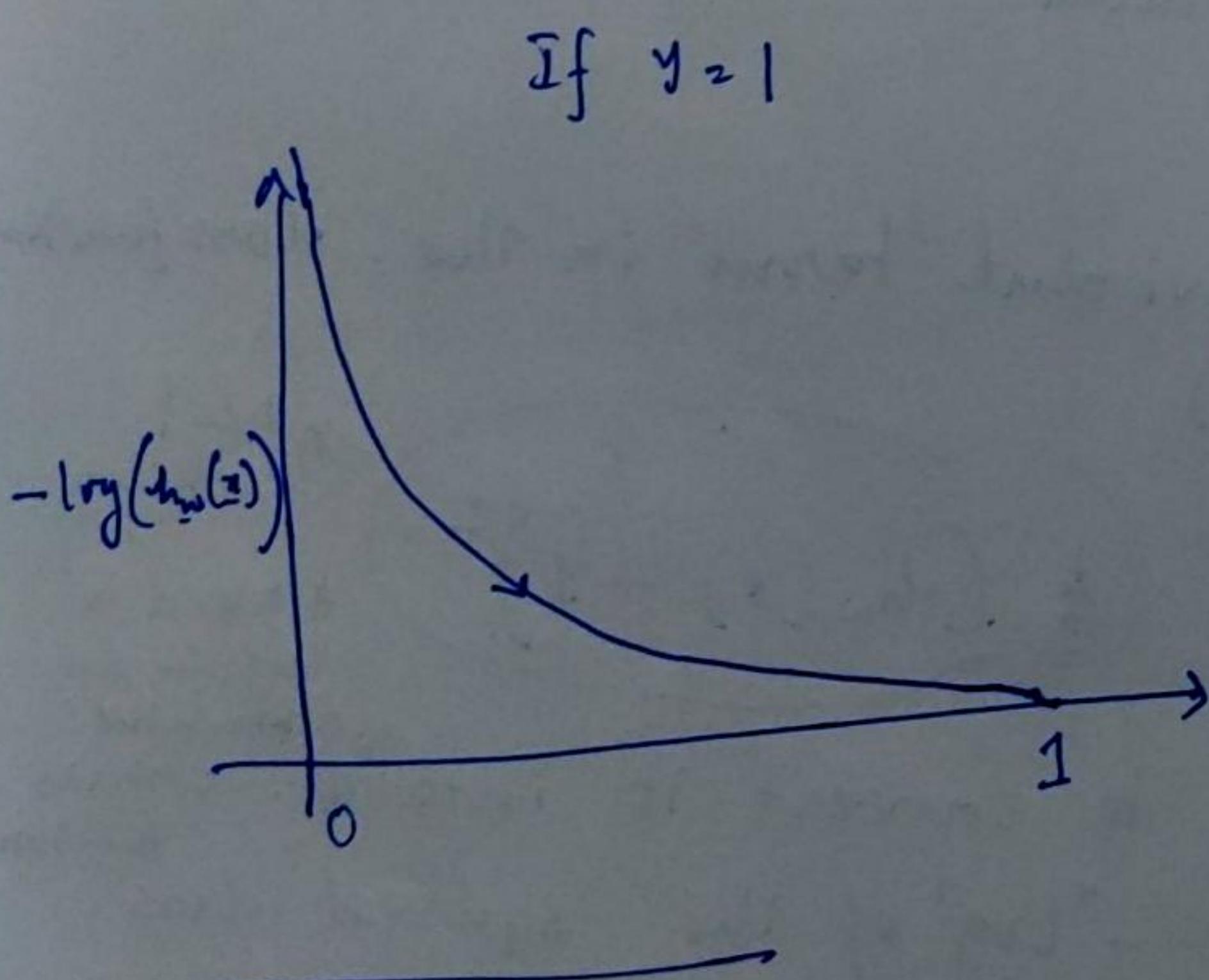
$$\boxed{\text{cost} = -\log \left(\frac{1-h_w(x)}{h_w(x)} \right)}$$

Case I: $h_w(x) \rightarrow 1$

$$-\log(1-h_w(x)) \rightarrow \infty$$

Case II: $h_w(x) \rightarrow 0$

$$-\log(1-h_w(x)) \rightarrow 0$$



Case I $h_w(x) = 1$

$$-\log(h_w(x)) = 0$$

Case II $h_w(x) \rightarrow 0$

$$-\log(h_w(x)) \rightarrow \infty$$

We can combine the two for binary classification cost function for binary classification.

$$\boxed{\text{cost}(h_w(x) - y) = -y \log(h_w(x)) - (1-y) \log(1-h_w(x))}$$

This gives the same cost for $y=0$ & $y=1$ as earlier.

Convex
function.

logistic regression : cost function

→ Binary classification

$$y^{(i)} = h_{\underline{w}}(\underline{x}^{(i)})$$

$$J(\underline{w}) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\underline{w}}(\underline{x}^{(i)}) - y^{(i)})$$

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\underline{w}}(\underline{x}^{(i)})) + (1-y^{(i)}) \log(1-h_{\underline{w}}(\underline{x}^{(i)})) \right]$$

**

The optimization prob. then boil down to.

minimize $J(\underline{w})$ or $\underline{w}^* = \underset{\underline{w}}{\operatorname{argmin}} (J(\underline{w}))$.

Please note $J: \mathbb{R}^n \rightarrow \mathbb{R}$.

output $h_{\underline{w}}(\underline{x}) = \left(\frac{1}{1+e^{-\underline{w}^T \underline{x}}} \right)$ w $\underbrace{P(y=c_k | \underline{x}; \underline{w})}_{k=2, \dots, n}$

Algorithms for optimization:

Type I: Gradient descent. (Batch).

$$\underline{w}^{k+1} = \underline{w}^k - \alpha \underline{d}^k$$

$$\text{w/ } \underline{d}^k = -\nabla_{\underline{w}} J(\underline{w})$$

$$\begin{aligned} &= - \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_j} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{pmatrix} = - \frac{\partial J}{\partial w_j} \quad \forall j = 1, 2, \dots, n \\ &\quad = \frac{1}{m} \sum_{i=1}^m (h_{\underline{w}}(\underline{x}^{(i)}) - y^{(i)}) \underline{x}_j^{(i)} \end{aligned}$$

Hence Loop over K
for i = 1, m

$$\underline{w}^{k+1} = \underline{w}^k - \alpha \left[\frac{1}{m} \sum_{i=1}^m (y^i - h_w(\underline{x}^i)) \underline{x}^i \right]$$

does this change anything?

$$y^i - h_w(\underline{x}^i)$$

$$(h_w(\underline{x}^i) - y^i)$$

Total sum over data set

Stop until converges.

(Stochastic gradient descent.)

Loop over k

for i = 1, m (loop over data).

$$\underline{w}^{k+1} = \underline{w}^k - \alpha \left(\frac{y^i - h_w(\underline{x}^i)}{m} \right) \underline{x}^i$$

end.

Stop until converges.

If not written in vectorized form \underline{w} etc.

Loop over k

for i = 1, m (loop over data)

for j = 1, n (loop over # of features) + 1

$$w_j^{k+1} = w_j^k - \alpha (y^i - h_w(\underline{x}^i)) x_j^i$$

end

stop until converges.

Other optimization algorithms.

- Discuss .
- i) Gradient descent
 - ii) Conjugate gradient
 - iii) BFGS
 - iv) L-BFGS.

Derivation of cost function (Max Likelihood approach)

** Another look of deriving the cost function from "maximum likelihood approach"

lets assume: $p(y=1 | \underline{x}; w) = h_w(\underline{x})$

$$p(y=0 | \underline{x}; w) = 1 - h_w(\underline{x})$$

we can assume a probability distribution.

for posteriors

$$p(y | \underline{x}; w) = (h_w(\underline{x}))^y (1 - h_w(\underline{x}))^{1-y}$$

Now, from this probability distribution the likelihood function is defined as.

$$L(w) = p(\underline{y} | \underline{X}; w)$$

$$= \prod_{i=1}^m p(y_i | \underline{x}_i; w) \quad \begin{array}{l} \text{[assuming} \\ \text{all the data} \\ \text{pts are indep]} \end{array}$$

$$= \prod_{i=1}^m (h_w(\underline{x}_i))^y (1 - h_w(\underline{x}_i))^{1-y}$$

To maximize the likelihood. (or log of the likelihood)

$$\log(L(w)) = \sum_{i=1}^m y^i \log(h_w(\underline{x}_i)) + (1-y^i) \log(1 - h_w(\underline{x}_i))$$

or minimize the $\frac{-\log(L(\underline{w}))}{m}$

Our optimization problem is posed as.

$$\begin{aligned} \min_{\underline{w}} J(\underline{w}) &= -\min_{\underline{w}} -\log(L(\underline{w})) \\ &= -\sum_{i=1}^m y^{(i)} \log h_{\underline{w}}(\underline{x}^{(i)}) + (1-y^{(i)}) \log(1-h_{\underline{w}}(\underline{x}^{(i)})) \end{aligned}$$

Now, to find $\nabla_{\underline{w}} J(\underline{w})$ we need. (for i)

$$\begin{aligned} \frac{\partial J}{\partial w_j} &= -\frac{\partial}{\partial w_j} \left[y^{(i)} \log g(\underline{w}^T \underline{x}^{(i)}) + (1-y^{(i)}) \log(1-g(\underline{w}^T \underline{x}^{(i)})) \right] \\ &= -\left[y^{(i)} \frac{1}{g(\underline{w}^T \underline{x}^{(i)})} g'(\underline{w}^T \underline{x}^{(i)}) + (1-y^{(i)}) \frac{1}{1-g(\underline{w}^T \underline{x}^{(i)})} (-g'(\underline{w}^T \underline{x}^{(i)})) \right] \end{aligned}$$

now we know $g'(z) = g(z)(1-g(z))$

Substituting we get

$$\begin{aligned} \frac{\partial J}{\partial w_j} &= [y^{(i)}(1-g(\underline{w}^T \underline{x}^{(i)})) - (1-y^{(i)})g(\underline{w}^T \underline{x}^{(i)})] x_j \\ &= (y^{(i)} - h_{\underline{w}}(\underline{x}^{(i)})) x_j \quad \forall j = 1, \dots, m \end{aligned}$$

The update rule turns out to be same
as in "regression prob"

$h_w(\underline{x}^{(i)})$ is now a nonlinear function
of $\underline{w}^T \underline{x}^i$

GLM: generalized linear models.
discusses the models which are linear
with respect to the "features" x_j

Note: If we combine all $j = 1, 2, \dots, m$ and vectorize.
the expression

$$\begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_m} \end{pmatrix} = (y^{(i)} - h_w(\underline{x}^{(i)})) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}^i$$

$$\Rightarrow \boxed{\nabla_w J = (y^{(i)} - h_w(\underline{x}^{(i)})) \underline{x}^{(i)}}$$

i is the data point.

Perceptron learning Algorithm

~ Perceptron algorithm has some historical significance.

~ The logistic regression method can be modified to change the nonlinear mapping

$$g(\underline{\omega}^T \underline{x} + \omega_0) \text{ as}$$

$$g(\underline{\omega}^T \underline{x} + \omega_0) = g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

~ If we then let $h_{\omega}(\underline{x}) = g(\underline{\omega}^T \underline{x} + \omega_0)$ as before, but using this modified definition of g , we can use the update rule.

$$\underline{\omega}_j \leftarrow \underline{\omega}_j + \alpha (y^{(i)} - h_{\omega}(\underline{x}^{(i)})) \underline{x}_j^{(i)}$$

This is "perceptron learning algorithm".

Note: In 1960, this perceptron was argued for how individual neurons in the brain work. Inspite. of its simplicity, "perceptron algorithm" is difficult to be endowed a "meaningful probabilistic interpretation" OR derive "perceptron" as a "maximum likelihood estimation algo".

Logistic Regression

"Multiclass classification"

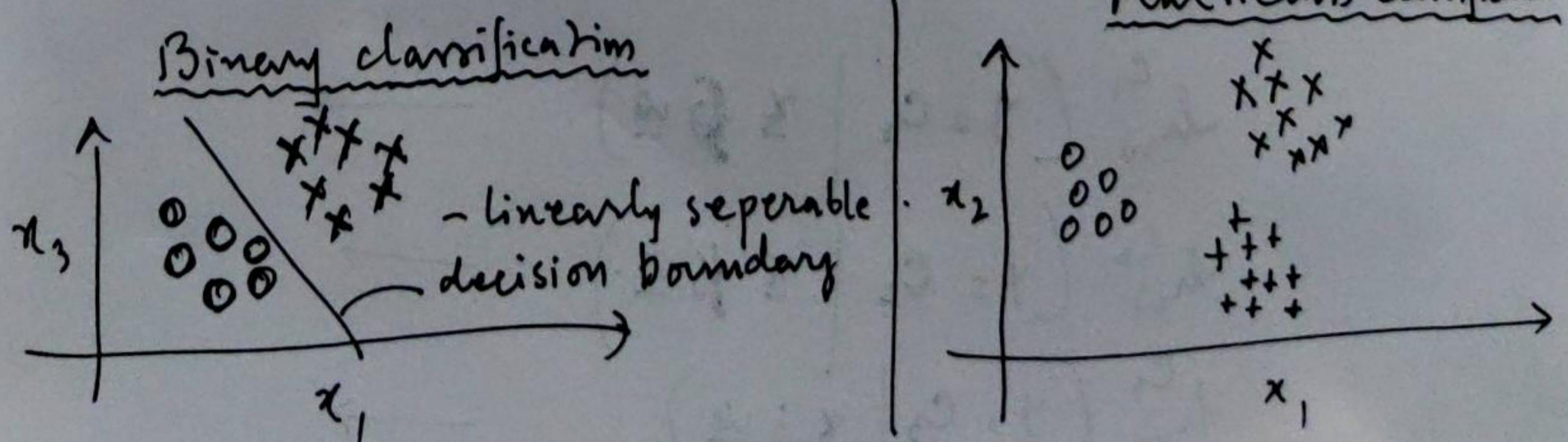
"One-vs-all" / On-vs-rest. (Ref Andrew Ng).

Email foldering/tagging: Work, family, Friends, Hobby.
 $y = c_1$ $y = c_2$ $y = c_3$ $y = c_4$

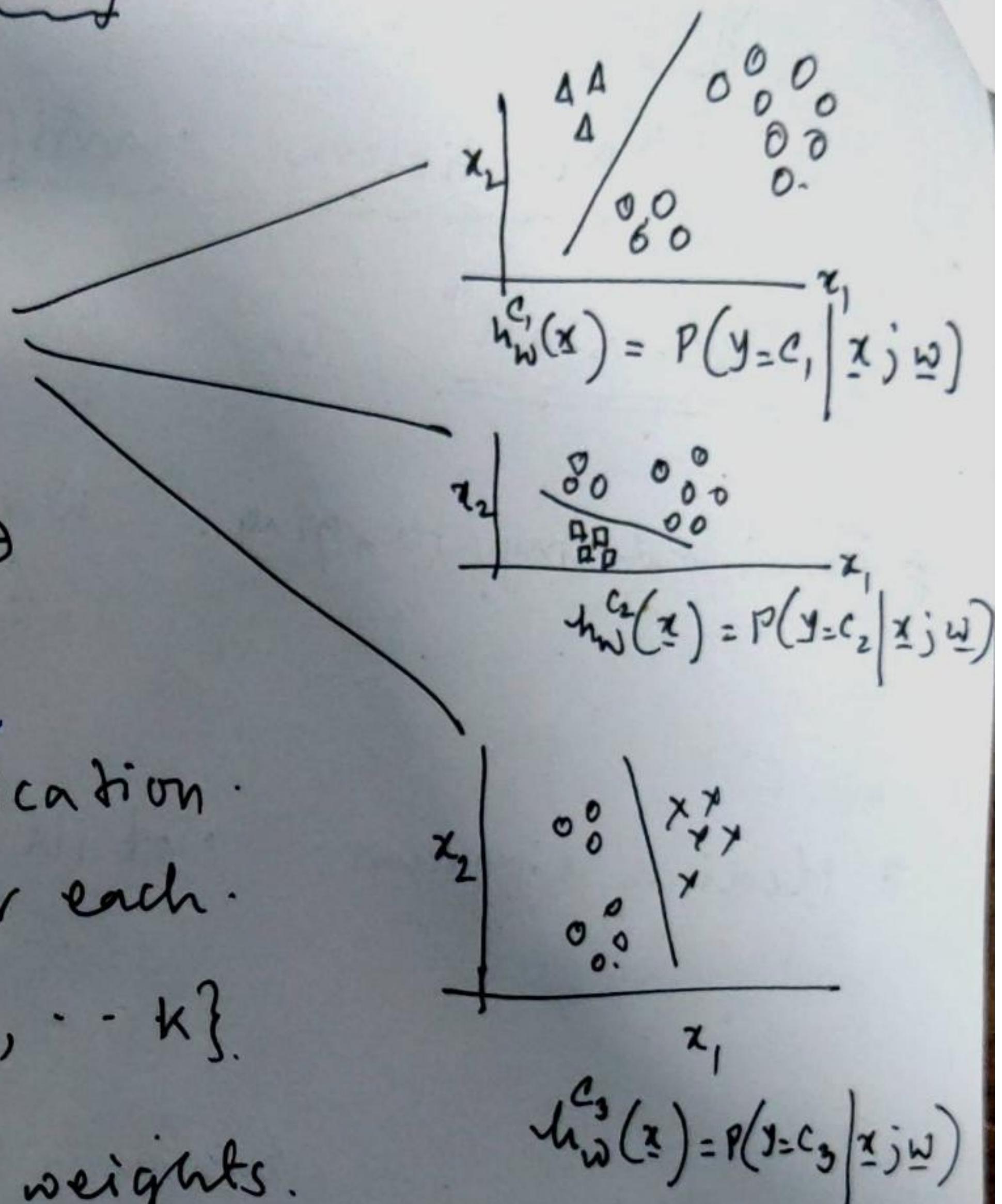
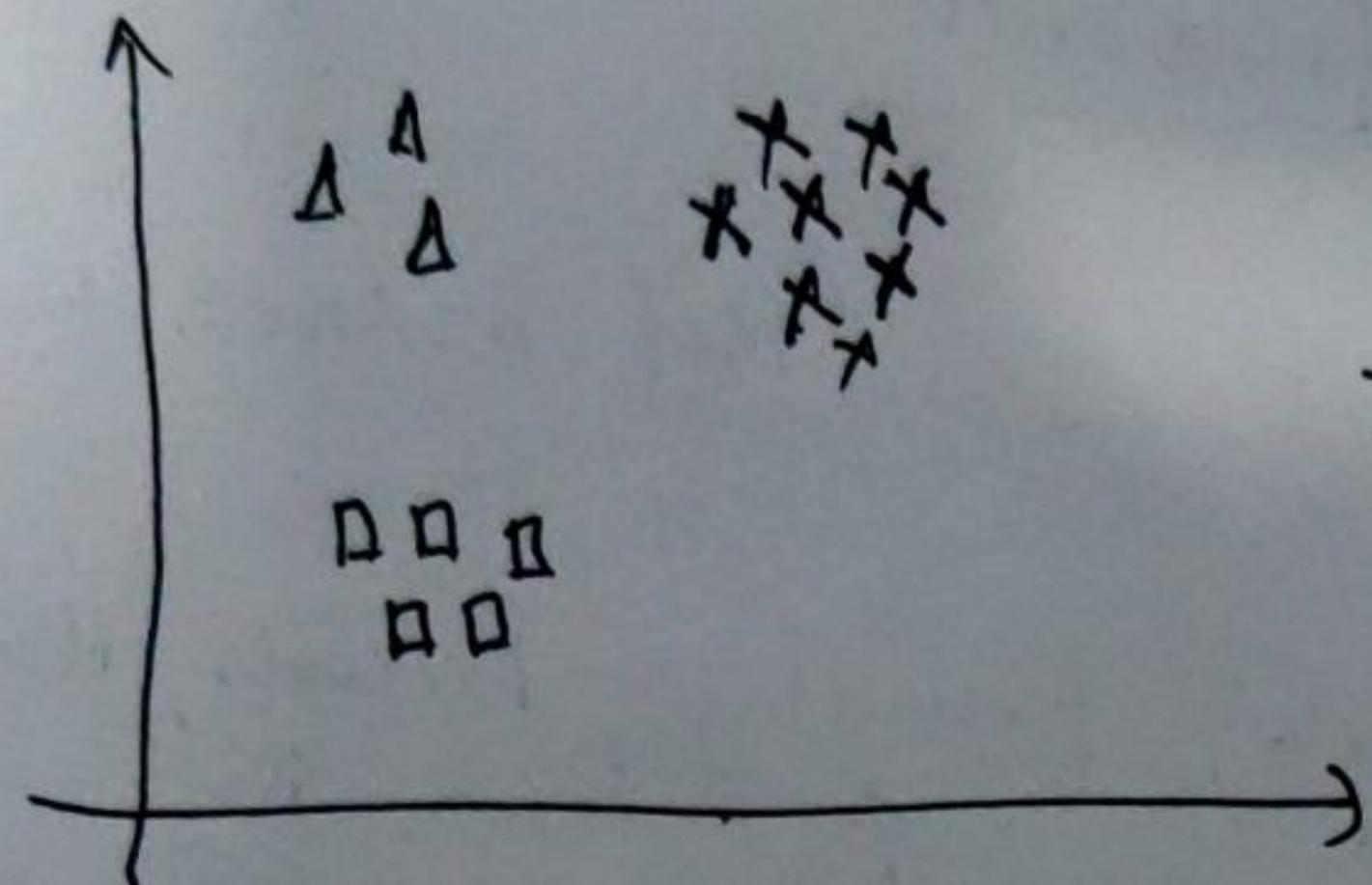
Medical diagram: Not ill, cold, Flu
 $y = c_1$, $y = c_2$, $y = c_3$

Weather: Sunny, Cloudy, Rain, Snow.
 $y = c_1$, $y = c_2$, $y = c_3$, $y = c_4$

When "data" can be labeled into multiple independent class, we need multi-class classification algorithm.



One-vs-all (One-vs-rest)



During the training phase

- Define binary classification probability density for each class C_k , $k \in \{0, 1, 2, \dots, K\}$.
- Find the optimized weights to define probability class conditional prob posterior prob

$$h_w^{C_k}(\underline{x}) = P(y = C_k | \underline{x}; \underline{w}) = \text{softmax}$$

Training

- So if there are 3 classes we have

$$h_w^{C_1}(y = C_1 | \underline{x}; \underline{w}) \rightarrow$$

$$h_w^{C_2}(y = C_2 | \underline{x}; \underline{w}) \rightarrow$$

$$h_w^{C_3}(y = C_3 | \underline{x}; \underline{w}) \rightarrow$$

Prediction

For prediction purpose, for a new input feature, \underline{x} , pick the class C_k , that maximize

$$\boxed{\max_{C_k} h_w^{C_k}(\underline{x})}$$

Summary: ~Binary classification applied one at a time to different classes.

~Training set is ~~arranging~~ arranged in such a way to. train the posterior prob. function. of the classes.