

### g. Cost function for logistic regression.

#### I. Training set for our Model.

Let  $m = \text{no. of training examples. } i = 1 \dots m$

$n = \text{no. of features } j = 1 \dots n$ .

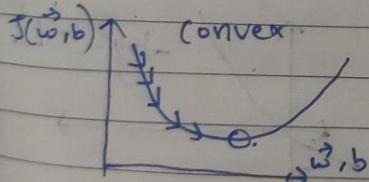
target.  $y$  is 0 or 1.

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

For linear regression, our squared error cost function looks like this:

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

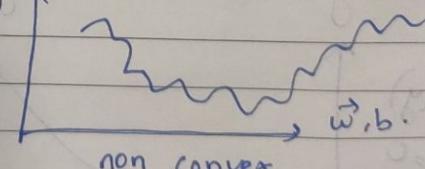
In case linear regression,  $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$ .



$$\text{for logistic regression, } f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

If we use gradient descent, it will have a lot of local minimas.

It won't guarantee global minima.



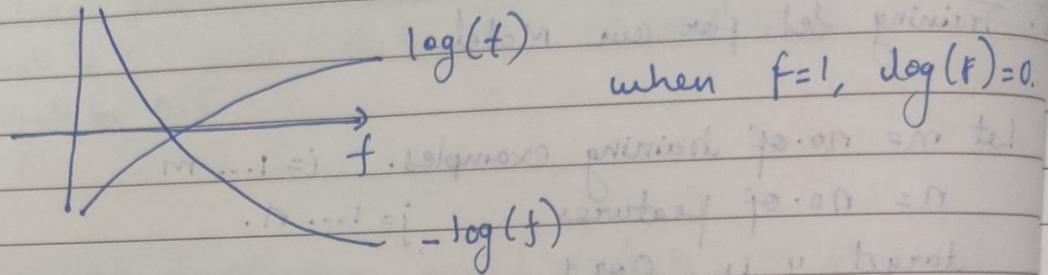
Let loss  $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$ . it is  $+^n$  of model  $f(\vec{x})$  and  $y$  true label.

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

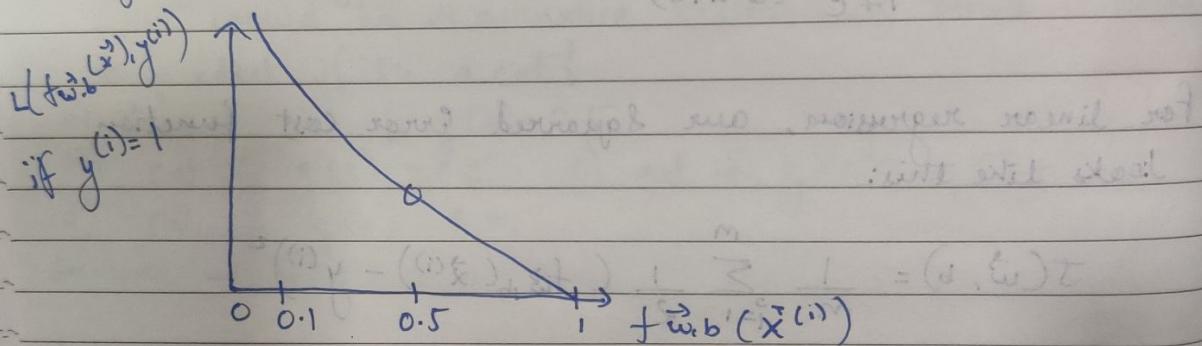
Logistic loss function.

Why does this make sense? (Graph of sigmoid function)

Now,



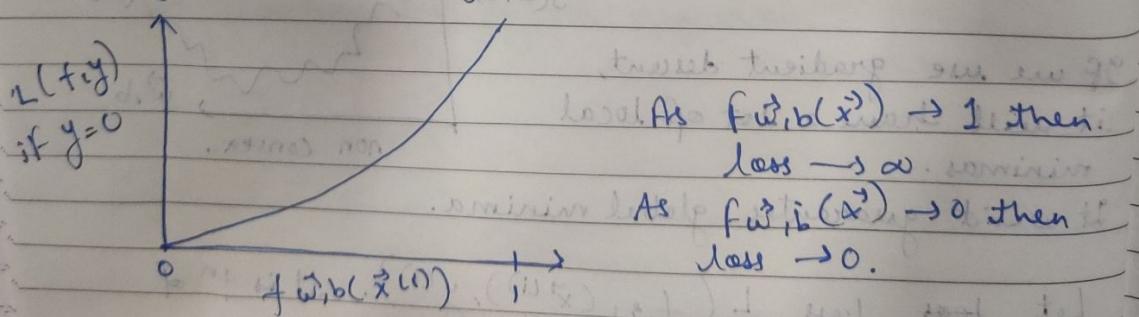
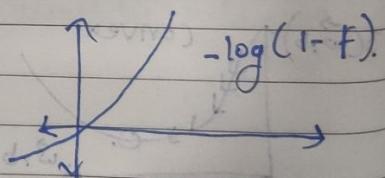
so,  $f$  is always between 0 and 1 for logistic regression.



As  $f_{w,b}(x^{(i)}) \rightarrow 1$ , then loss  $\rightarrow 0$ .

As  $f_{w,b}(x^{(i)}) \rightarrow 0$  then loss  $\rightarrow \infty$

when  $y^{(i)} = 0$ , our log function is



the further prediction  $f_{w,b}(x^*)$  is from the target  $y^{(i)}$ , the higher the loss.

## II. Cost.

$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}, y^{(i)}))$$

loss.

when  $L = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1, \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0. \end{cases}$

Using this, our cost function can reach a global minimum.

Now, find  $\vec{w}, b$  that minimize  $J$ .

Question: Why is the squared error cost not used in logistic regression?

- The non-linear nature of the model results in a "wiggly", non convex cost function with many potential local minima.
- The mean squared error is used for logistic regression.

Answer- a.

## III. Simplified Cost function.

As we know, our loss function looks like above.

Now,

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)})).$$

This equation will be our simplified ~~cost~~<sup>loss</sup>, since  $y$  can only take values of 1 and 0.

if  $y^{(i)} = 1$

$$L(f_{\vec{w}, b}(\vec{x}), y^{(i)}) = -\log(f(\vec{x})). \quad \left. \right\} \text{loss values.}$$

if  $y^{(i)} = 0$

$$L(f(\vec{x}), y) = -\log(1-f(\vec{x})). \quad \left. \right\} \text{marked}$$

The simplified cost function will be.

$$\begin{aligned} J(\vec{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})] \\ &= -\frac{1}{m} \left[ y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)})) \right] \end{aligned}$$

This cost function is derived from Maximum Likelihood  
Parzen window statistics.

The cost function above is a convex function.

Question: The above simplified cost function, if the target  $y^{(i)} = 1$ , then what does the expression simplify to?

- a)  $-\log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))$       b)  $-\log(f_{\vec{w}, b}(\vec{x}^{(i)}))$

Answer - b.