

LOGISTIC REGRESSION

Classification

- Email: Spam / Not Spam?
- Online Transactions: Fraudulent (Yes/No)?
- Tumor: Malignant / Benign.

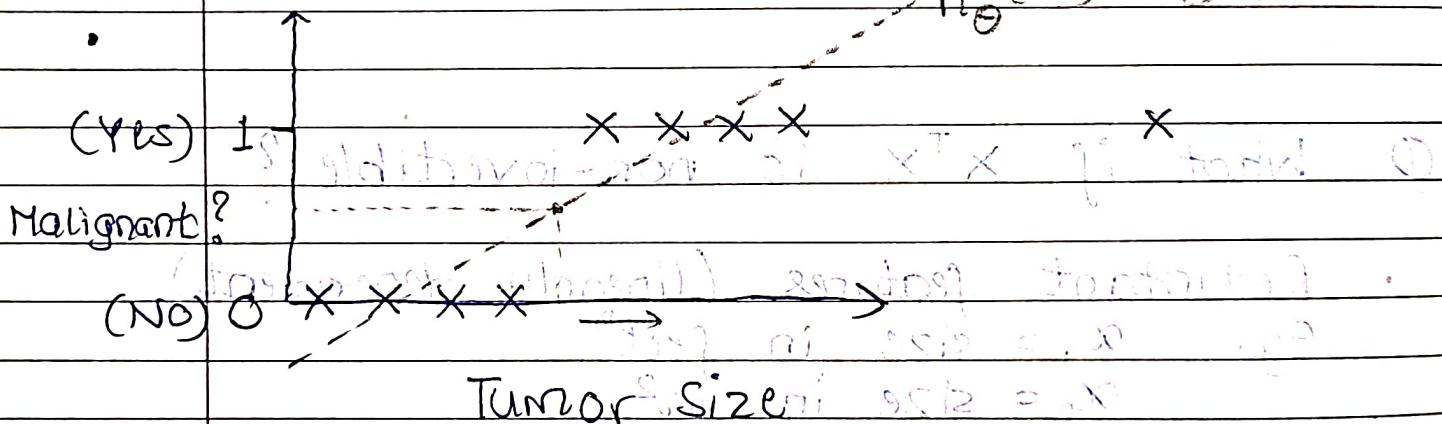
$y \in \{0, 1\}$ 0: "Negative Class" (Benign)

1: "Positive Class" (Malignant)

y may even take more values than 0 or 1

$y \in \{0, 1, 2, 3\} \rightarrow$ Multiclass Classification.

$$h_{\theta}(x) = \theta^T x$$



Now, one thing that we could do is to apply Linear Regression to this data set and try to fit a straight line.

$$h_{\theta}(x) = \theta^T x$$



If you want to make predictions from the line fit to the data set, you can set a threshold where the classifier outputs 1 for values greater than 0.5 (vertical axis value) and if the hypothesis outputs a value less than 0.5, then you can take $y=0$.

Threshold classifier output $h_0(x)$ at 0.5:

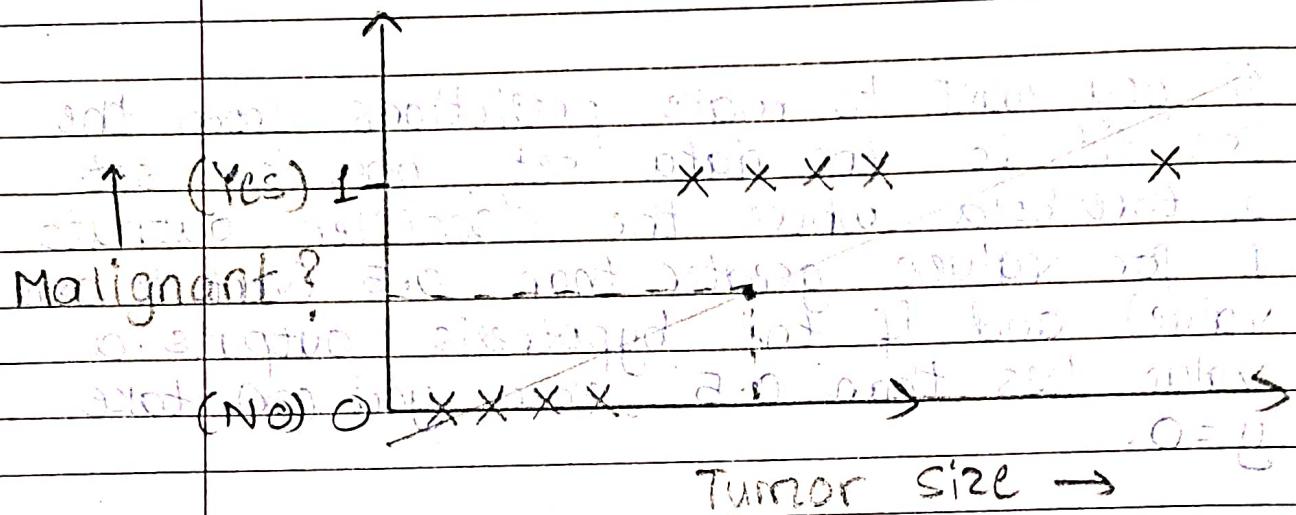
if $h_0(x) \geq 0.5$, predict $y=1$

else if $h_0(x) < 0.5$, predict $y=0$

Till here, it looks like linear regression is giving us some useful results even in a classification problem.

But now, let's add one more training example.

Notice that, intuitively the new training example doesn't change anything and it's pretty clear that the tumor is malignant. But, now when we fit our line to the data, our model might come out to be somewhat like this



And now, clearly, this hypothesis is not accurate since it will predict a malignant tumor as \geq benign since it falls in the (< 0.5) category. So, linear regression fails here.

You might get lucky if you use it.

Also, $h_0(x)$ can be > 1 or < 0 .

Logistic Regression is a classification algorithm and we will use it.

- ii) Logistic Regression Model

Want $h_0(x) \leq 1$ and $h_0(x) \geq 0$

$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$

Logistic Regression to Model Malignant

$$h_{\theta}(x) = g(\theta^T x)$$

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

Sigmoid / Logistic Function

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

\therefore The hypothesis is a function of $\theta^T x$ that is $g(\theta^T x)$ where $g(z) = \frac{1}{1 + e^{-z}}$

Clearly, $g(z)$ will be in range [0,1] for all values of z from $-\infty$ to $+\infty$.

- Interpretation of Hypothesis Output.

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x .

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$$h_{\theta}(x) = 0.7 \quad \text{and } 70\%$$

Tell patient that 70% chance of tumor being malignant.

$$h_{\theta}(x) = P(y=1 | x; \theta)$$

"probability" that $y=1$, given x , parameterized by θ ".

$$(x^T \theta) p = (x^T \theta)_0$$

Now, the problem is, y can only be either 0 or 1

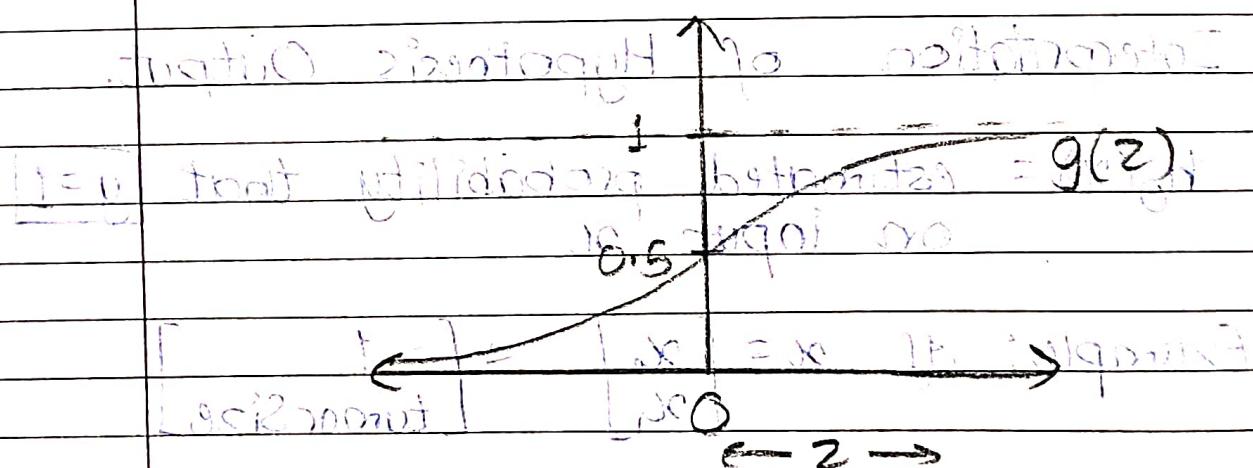
$$P(y=0|x; \theta) = 1 - P(y=1|x; \theta)$$

$$1 - e^{(x^T \theta)}$$

Decision Boundary

$$(x^T \theta) h_{\theta}(x) = g(\theta^T x) = \frac{1}{1+e^{-z}}$$

Suppose, prediction "y=1" if $h_{\theta}(x) > 0.5$
predict "y=0" if $h_{\theta}(x) < 0.5$



$$g(z) \geq 0.5 \quad \text{if } z = (x)^T \theta$$

when $z \geq 0$

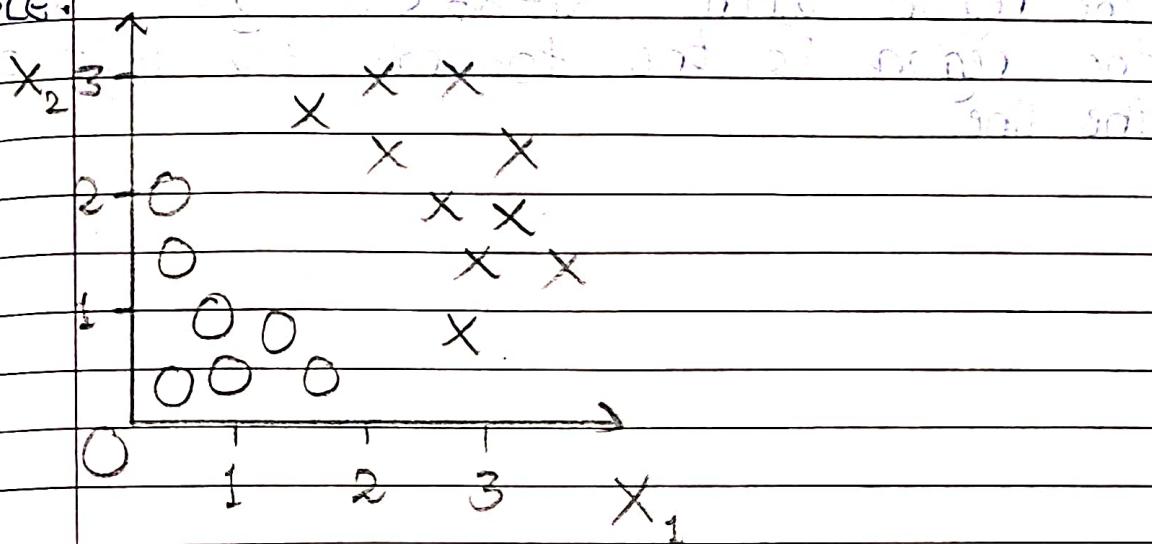
$$\therefore h_{\theta}(x) = g(\theta^T x) \geq 0.5 \text{ when } \theta^T x \geq 0$$

$$\text{whenever } \theta^T x \geq 0 \Rightarrow p = g(\theta^T x)$$

and $h_{\theta}(x) \leq g(\theta^T x) < 0.5$

whenever $\theta^T x < 0$

EXAMPLE:



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Say we choose $\theta_0 = -3$ $\theta_1 = 1$ $\theta_2 = 1$

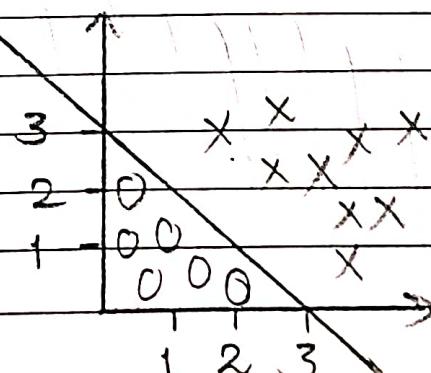
$$\Rightarrow \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Predict "y=1" if $\underbrace{-3 + x_1 + x_2}_{(\theta^T x)} \geq 0$

$$\Rightarrow x_1 + x_2 \geq 3$$

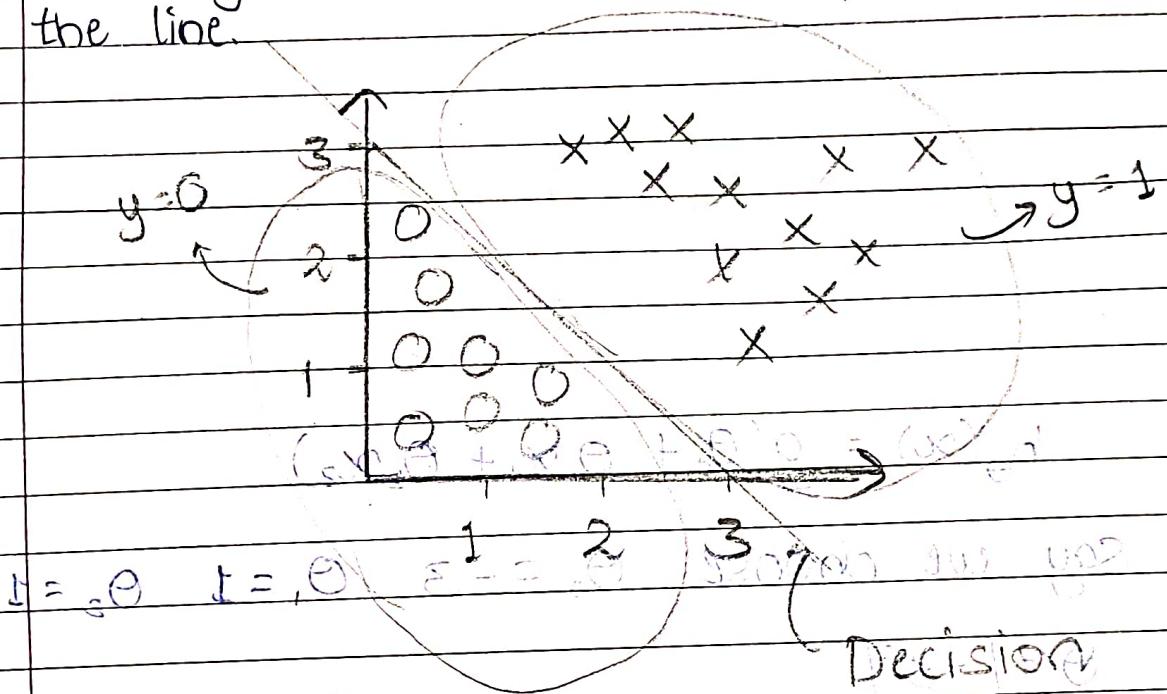
Let's see what that means on the figure

$$x_1 + x_2 = 3$$



∴ $x_1 + x_2 \geq 3$ means (the right) half
that is, everything to the up and right of
the line.

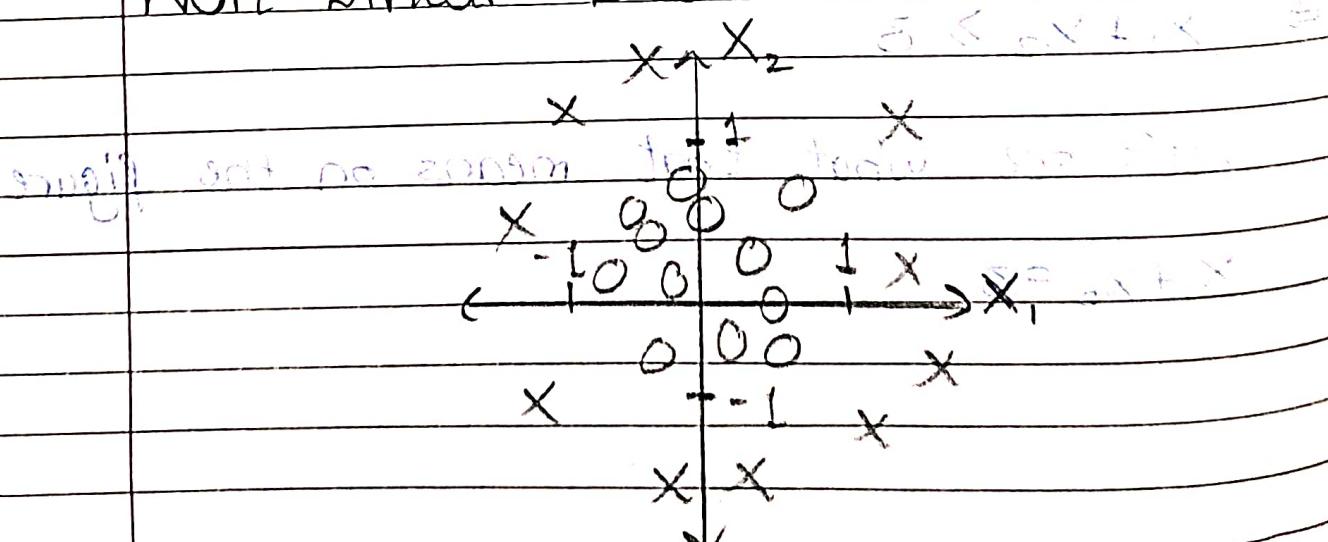
The region where $x_1 + x_2 < 3$, $y=0$, and
the region is the towards the left of
the line.



Decision
Boundary.

Here, we can also remove the data
set.

Non-Linear Decision Boundaries



Let's say our hypothesis looks like:

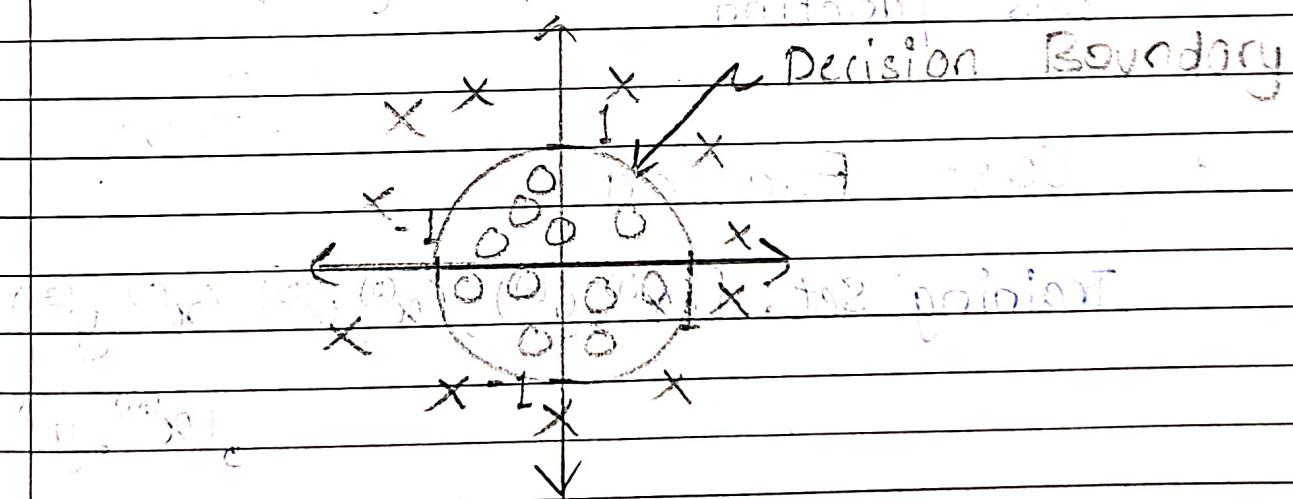
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

for now, assume that after some valid procedure to be specified, I end up choosing

$$\theta_0 = -1, \theta_1 = 0, \theta_2 = 0, \theta_3 = 1, \theta_4 = 1$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

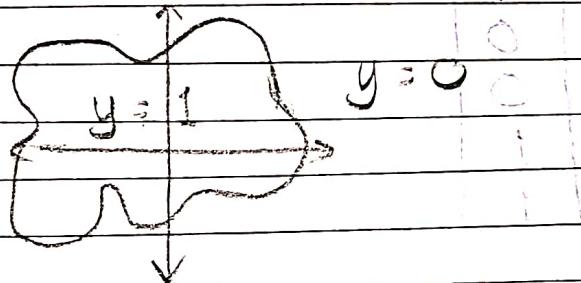
Predict "y=1" if $-1 + x_1^2 + x_2^2 \geq 0$



This will mean, everything outside the circle is the "y=1" region, and everything inside is "y=0" region.

NOTE: The decision boundary is a property, not of the training set, but the hypothesis under ~~(S)~~ the parameters. Thus, our vector θ will define the decision boundary.

This is amazing! Since we can get almost any type of decision boundary like this.



And this complexity increases proportionally with the complexity of our hypothesis function.

• Cost Function

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

train some examples in index	$x_{(0)}$	$x_{(1)}$	$x_{(2)}$	$x_{(3)}$	$x_{(4)}$	$x_{(5)}$	$x_{(6)}$	$x_{(7)}$	$x_{(8)}$	$x_{(9)}$
points vs label	0, 0, 0	1, 0, 0	0, 1, 0	0, 0, 1	1, 1, 0	1, 0, 1	0, 1, 1	1, 1, 1	0, 0, 0	1, 1, 1
	0, 0, 0	1, 0, 0	0, 1, 0	0, 0, 1	1, 1, 0	1, 0, 1	0, 1, 1	1, 1, 1	0, 0, 0	1, 1, 1
	0, 0, 0	1, 0, 0	0, 1, 0	0, 0, 1	1, 1, 0	1, 0, 1	0, 1, 1	1, 1, 1	0, 0, 0	1, 1, 1
	0, 0, 0	1, 0, 0	0, 1, 0	0, 0, 1	1, 1, 0	1, 0, 1	0, 1, 1	1, 1, 1	0, 0, 0	1, 1, 1

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameter?

In linear regression,

$$\text{Cost Function } J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \left\{ h_{\theta}(x^{(i)}) - y^{(i)} \right\}^2$$

$$\text{Let cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\therefore J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

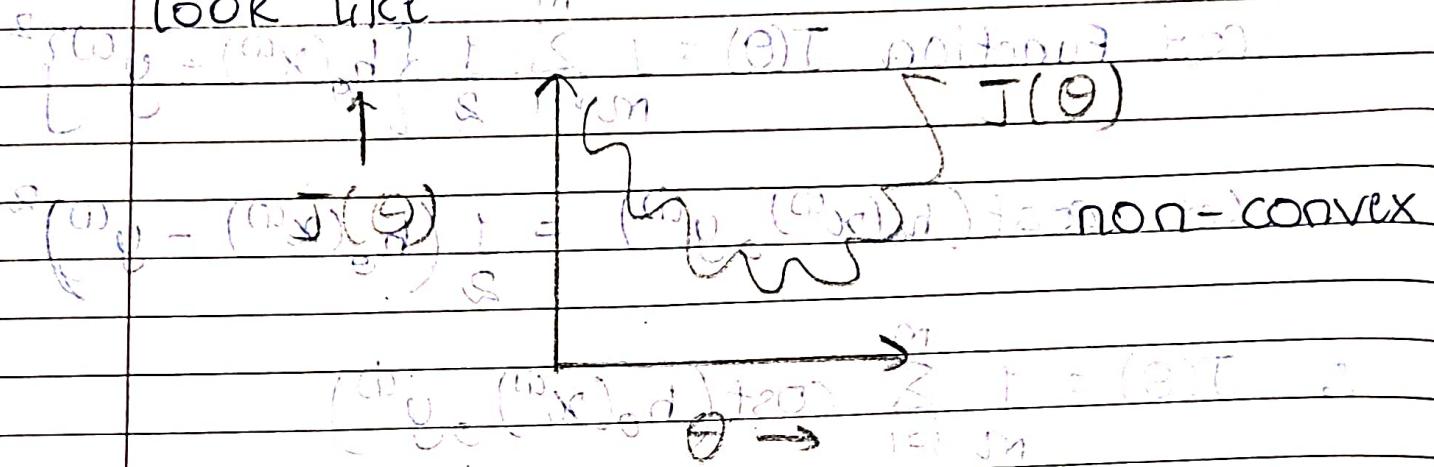
This is the cost the learning algorithm will have to pay if it outputs $h_{\theta}(x)$ where the correct value was y .

Now, this cost function works well for linear regression, but here, we're interested in logistic regression.

If we minimize of $\text{cost}(h(x), y)$ and then plug it in the $J(\theta)$ then linear regression will work ok. But, it turns out that the $h_{\theta}(x)$ function has non-linearity which can be easily seen by the fact that

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

So, the non-linear sigmoid function is plugged in the cost function and then is the $J(\theta)$ function and then plot to see what $J(\theta)$ looks like, we find it can look like



non-convex and it's easy to see that if we run gradient descent at a given point, it's not guaranteed that we will reach the global minima.

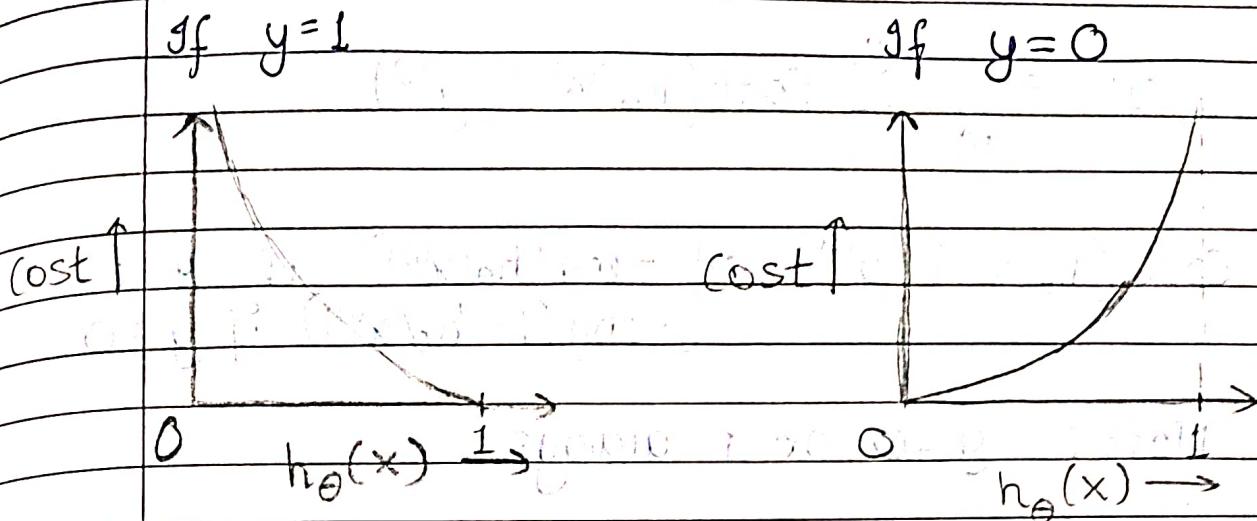
But we want a convex $J(\theta)$

(e) A zygote is the egg of most plants.

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

$$y_0 = \begin{cases} 1 & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

00122503 215201 21



$\text{cost}(y, \hat{y}) = (y - \hat{y})^2$ or $= (\hat{y} - y)^2$

If $y = 1$, $h(x) = 1 \rightarrow \text{cost} = 0$

If $y = 1$, $h(x) = 0 \rightarrow \text{cost} \rightarrow \infty$

If $y = 0$, $h(x) = 1 \rightarrow \text{cost} \rightarrow \infty$

This means, if the algorithm predicts certainly that the output is 0 or the $h_{\theta}(x) = 0$ but if actually $y = 1$ or the actual output is 1, then we 'penalize' the model by $\text{cost} = \infty$ from the graph (very large cost).

Similarly, if $y = 0$ & $h_{\theta}(x) = 1$, the cost will be extremely large $\rightarrow \infty$ and we will penalize (the) algo by a very large value.

$$\begin{aligned} \text{Cost}(h_{\theta}(x), y) &= 0 \quad \text{if } h_{\theta}(x) = y \\ \text{Cost}(h_{\theta}(x), y) &\rightarrow \infty \quad \text{if } \begin{cases} y = 0 \text{ & } h_{\theta}(x) = 1 \\ y = 1 \text{ & } h_{\theta}(x) = 0 \end{cases} \end{aligned}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_\theta(x^{(i)})) & \text{if } y = 1 \\ -\log(1 - h_\theta(x^{(i)})) & \text{if } y = 0 \end{cases}$$

Note: $y = 0$ or 1 always

$$\text{Cost}(h(x), y) = -y \log(h(x)) - (1-y) \log(1-h(x))$$

$$\text{If } y = 1, \text{ cost}(h(x), y) = -\log(h(x))$$

$$\text{If } y = 0, \text{ cost}(h(x), y) = -\log(1-h(x))$$

$$\text{So } J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$J(\theta) = -\frac{1}{m} \left\{ \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right\}$$

$$J(\theta) = \frac{1}{m} \left\{ \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right\}$$

This is nice as this is convex

To fit parameters θ :

$$\min_{\theta} J(\theta)$$

Get θ from iteration

To make a prediction given new x_i

$$\text{Output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} (P(y=1|x; \theta))$$

GRADIENT DESCENT:

$$\text{Want } \min_{\theta} J(\theta)$$

Repeat { until (all θ_j updated)}

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

(simultaneously update all θ_j)

where $J(\theta)$ is the function defined before

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(simultaneously update all θ_j)

The difference is $h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$

Vectorized Implementation:

$$h = g(X\theta)$$

where $X \in \mathbb{R}^{n \times m}$
 $\theta \in \mathbb{R}^{m \times 1}$

$$J(\theta) = \frac{1}{m} \cdot \left(-y^T \log(h) - (1-y)^T \log(1-h) \right)$$

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - y)$$

Advanced Optimization Algorithms

Cost function $J(\theta)$. Want $\min_{\theta} J(\theta)$

Given θ , we have code that can compute

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$ (for $j = 0, 1, \dots, n$)

Gradient Descent:

Repeat

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

(simultaneous)

So, for gradient descend, we need code to supply $J(\theta)$ {for monitoring the correctness of g.d.} and $\frac{\partial}{\partial \theta_j} J(\theta)$

and plug in the value in the g.d. function to get the updated values of θ for minimizing $J(\theta)$.

There are many more optimization algos.:

- Conjugate Gradient
- BFGS
- L-BFGS

Advantages :

- No need to manually pick α
- Often faster than g.d. and the

These algorithms have a clever inner-loop (like a line search algorithm) that automatically tries out different values of α and then chooses a good learning rate. So, it can even pick different values of α for every iteration automatically.

They also do much more than just picking a good learning rate α and so, they often end up converging much faster than gradient descent.

Disadvantage:

- Complex

Implementation without any library

Explanation with an Example:

$$\Theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\Theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial J(\Theta)}{\partial \theta_1} = 2(\theta_1 - 5)$$

$$\frac{\partial J(\Theta)}{\partial \theta_2} = 2(\theta_2 - 5)$$

Ans: $\theta_1 = 5$ $\theta_2 = 5$

But how to do it with some optimization algorithm.

```

jVal = 1 * (theta(1) - 5)^2 + 1 * (theta(2) - 5)^2;
gradient = zeros(2,1);
gradient(1) = 2 * (theta(1) - 5);
gradient(2) = 2 * (theta(2) - 5);

options = optimset('GradObj','on','MaxIter',10)
initialTheta = zeros(2,1);

```

$[optTheta, functionVal, exitFlag] =$

$fminunc(@costFunction, initialTheta, options)$

For the implementation of this in Octave,
the θ vector must be at least a 2-dimensional vector.

$$\theta \in \mathbb{R}^d \quad | \quad d \geq 2$$

$$\text{theta} = [\theta_0; \theta_1; \theta_2; \dots; \theta_{n+1}]$$

function [jVal, gradient] = costFunction(theta)

jVal = [code to compute $J(\theta)$];

gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

⋮

gradient($n+1$) = [code to compute $\frac{\partial}{\partial \theta_{n+1}} J(\theta)$];

- notes -

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

Multiclass Classification

Emails: Work, Friends, Family, Hobby

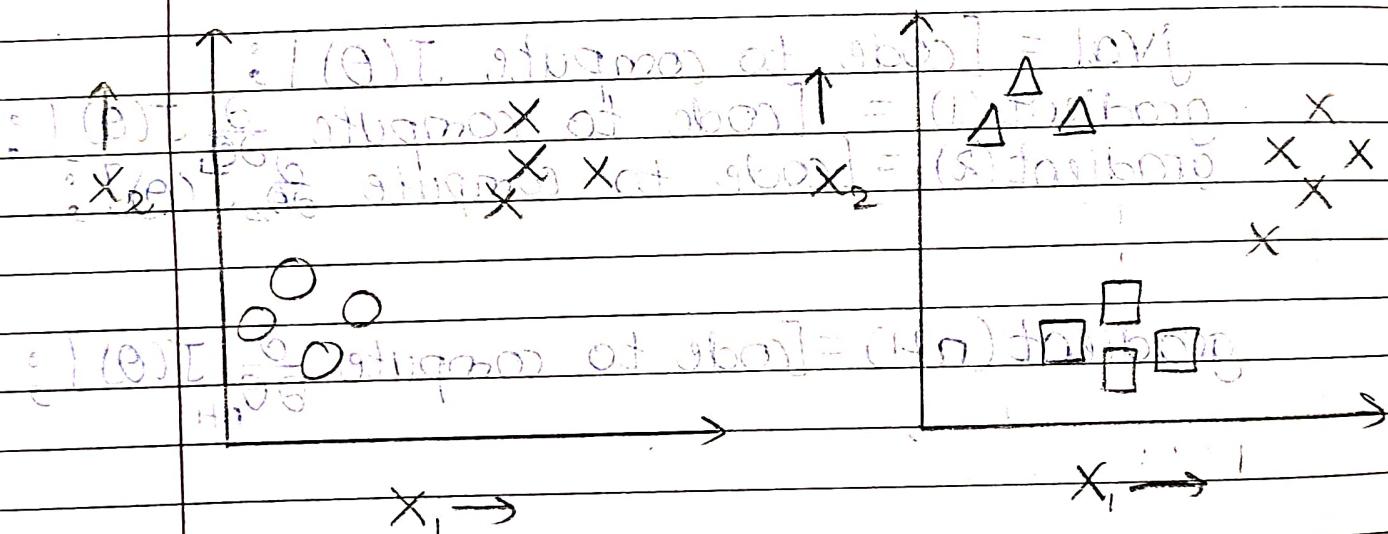
$y=1$ Work $y=2$ Friends $y=3$ Family $y=4$ Hobby

Medicine: Not ill, Cold, Flu, Headache

Weather: Sunny, Cloudy, Rainy, Snowy

Binary Classification: Multi-class Classification

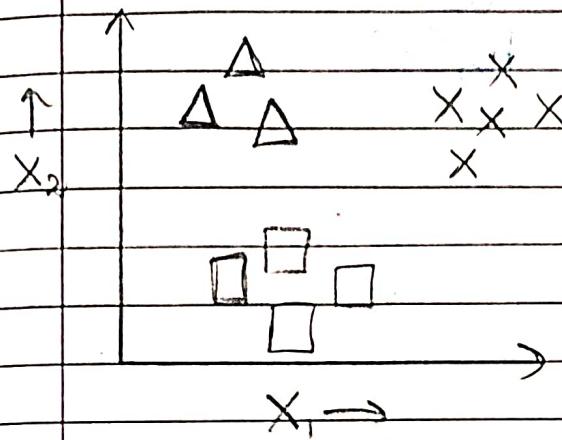
Intuition: Multiclass = [Binary, Binary, ...]



Solution -

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

ONE - VS - ALL (ONE - VS - REST) :



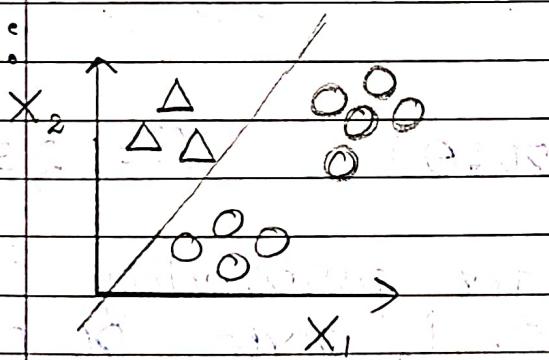
Class 1: Δ

Class 2: \square

Class 3: \times

We take our training set and turn it into 3 separate two class classification problems.

Class 1:

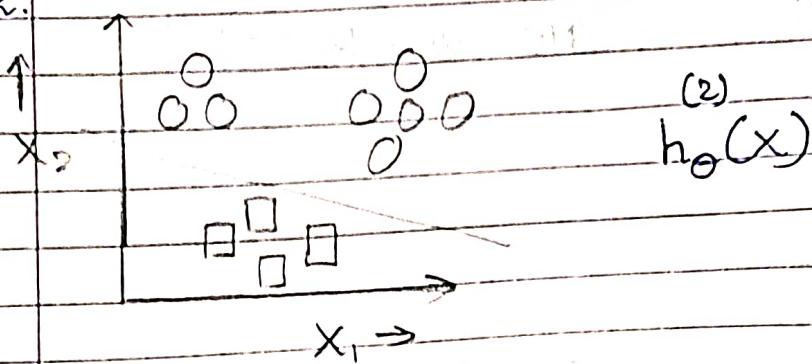


So, here this is a sort of fake training set for class 1 where class 1 is the positive class and classes 2 and 3 get assigned to the negative class.

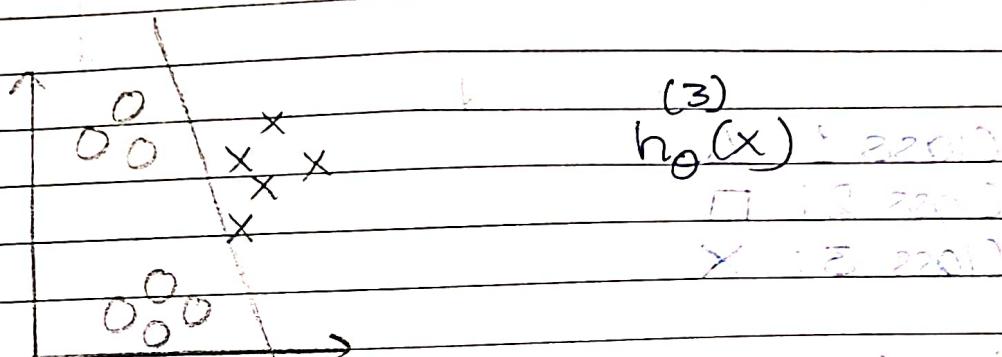
Now, we will fit a classifier, say $h_{\theta}^{(1)}(x)$

$$h_{\theta}^{(1)}(x)$$

Class 2:



Class 3:



We have fit 3 individual classifiers for the 3 classes. So, for $i=1, 2, 3$, we have fit a classifier

$$h_i(x) = P(y=i|x; \theta) \quad (i=1, 2, 3)$$

which estimates the probability that y is equal to class i , given x parameterised by θ .

And then, we will pick the value i that maximizes $h_i(x)$.

$$\max_i h_i(x)$$

Regularized Logistic Regression

We have previously discussed that even logistic regression is prone to overfitting by making the decision boundary have a large variance in order to make it fit all the existing training examples, resulting in failure to generalize the output for new data. This can happen when there are a large number of features and their higher-powered products in the $h(x)$.

Eg

$$\begin{aligned}
 & \text{Input Data} \\
 & \begin{array}{c|ccccc}
 & x_0 & x_1 & x_2 & x_3 & x_4 \\
 \hline
 1 & 0 & 0 & 0 & 0 & 0 \\
 2 & 0 & 1 & 0 & 0 & 0 \\
 3 & 0 & 0 & 1 & 0 & 0 \\
 4 & 0 & 0 & 0 & 1 & 0 \\
 5 & 0 & 0 & 0 & 0 & 1
 \end{array}
 \end{aligned}$$

$$\begin{aligned}
 & \text{Hypothesis Function: } h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots) \\
 & \text{Term 1: } \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \\
 & \text{Term 2: } \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \\
 & \text{Term 3: } \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \\
 & \text{Term 4: } \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \\
 & \text{Term 5: } \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots
 \end{aligned}$$

Cost Function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

After regularization, the cost function becomes:

$$\begin{aligned}
 J(\theta) = & -\left[\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))] \right] \\
 & + \frac{\lambda}{2m} \sum_{j=1}^n (\theta_j)^2
 \end{aligned}$$

Minimizing this cost function will lead to smaller values of θ , causing a simpler hypothesis.

Gradient Descent:

Repeat {

$$\theta_0 := \theta_{0,0} - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_{j,0} - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

($j = 1, 2, 3, \dots, n$)

}

Repeat {

$$\theta_0 := \theta_{0,0} - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_{j,0} \left(1 - \alpha \lambda \right) - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

}

MATH BEHIND LOGISTIC REGRESSION :

$$\hat{y} = \frac{1}{1 + e^{-(\theta^T x)}} \quad \text{where } \theta \in \mathbb{R}^{(n+1) \times 1} \text{ and } x \in \mathbb{R}^{m \times (n+1)}$$

Now, for logistic regression -

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

For simplicity, let's ignore m and let's vectorize this function.

$$J = -\{ \mathbf{y}^T \log(\mathbf{h}) + (1-\mathbf{y})^T \log(1-\mathbf{h}) \}$$

where I have darkened the vectors,

$$\mathbf{y} \in \mathbb{R}^{(m)}$$

$$\mathbf{h} \in \mathbb{R}^{(m)}$$

We want to find the derivative of this cost function with respect to θ (b and w)

$$\min_{\theta, w} J = \frac{\partial J}{\partial \theta}$$

Let's first find $\frac{\partial J}{\partial w}$

Clearly, by chain rule,

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial h} \cdot \frac{\partial h}{\partial w}$$

To make this clearer, assume that

$$① \quad a = h(\vec{y})$$

$$② \quad z = \vec{w}^T \vec{x} + b$$

$$\text{so, } a = \sigma(z) \Rightarrow \frac{e^z}{1+e^z} = a \Rightarrow e^z = \frac{a}{1-a}$$

$$\therefore \frac{\partial J}{\partial w} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\Rightarrow \frac{\partial J}{\partial a} = - \left\{ \frac{\vec{y}^T}{a} - \frac{(1-\vec{y})^T}{(1-a)} \right\} = \left[\frac{-\vec{y}^T}{a} + \frac{(1-\vec{y})^T}{(1-a)} \right]$$

$$\frac{\partial a}{\partial z} = \frac{\partial}{\partial z} \left\{ \frac{1}{1+e^{-z}} \right\} = -(1+e^{-z})^{-2} \cdot -e^{-z}$$

$$\frac{\partial a}{\partial z} = \frac{e^{-z}}{(1+e^{-z})^2} = (1-a) \times \frac{1}{a} = \frac{1}{a(1-a)}$$

$$\frac{\partial z}{\partial w} = \vec{x}$$

$$\frac{\partial J}{\partial w} = \left\{ \frac{-\vec{y}^T}{a} + \frac{(1-\vec{y})^T}{(1-a)} \right\} \{a(1-a)\} \{ \vec{x} \}$$

$$\frac{\partial J}{\partial w} = \left\{ \frac{-(1-a)\vec{y}^T}{a(1-a)} + \frac{(1-\vec{y})^T a}{a(1-a)} \right\} \{ \vec{x} \}$$

$$\frac{\partial J}{\partial w} = (-\vec{y}^T + \vec{y}^T a + a - \vec{y}^T a) \vec{x}$$

$$\boxed{\frac{\partial J}{\partial w} = (A - Y) \vec{x}} \approx \vec{x}^T \cdot (A - Y) \quad (n \times m) \quad (m \times 1)$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial b}$$

$$\frac{\partial J}{\partial b} = \left\{ -y^T + (1-y)^T \right\} \left\{ a(1-a) \right\} \left\{ 1 \right\}$$

$$\frac{\partial J}{\partial b} = -y^T(1-a) + (1-y)^T a = -y^T + y^T a + a - y^T$$

$$\frac{\partial J}{\partial b} = (A - Y) \approx \text{np.sum}(A - Y)$$

$$w = w - \alpha \cdot \frac{\partial J}{\partial w}$$

$$b = b - \alpha \frac{\partial J}{\partial b}$$

$$(p-1)b = 1 \times (p-1) = 9 = 56$$

$$(\bar{x}) = \frac{56}{106}$$

$$(\bar{x})(p-1)b = 1 \times (p-1) + (p-1) = 16$$

$$(\bar{x})(p-1)b = 1 \times (p-1) + (p-1) = 16$$

$$x(p-1) + (p-1) = 16$$

$$(Y - A) \cdot X = \underbrace{X(Y - A)}_{(m \times n) \times (n \times m)} = 16$$