

Machine Learning for Engineering and Science Applications
Professor Dr. Balaji Srinivasan
Department of Mechanical Engineering
Indian Institute of Technology, Madras
Gradient of Logistic Regression

(Refer Slide Time: 0:25)

The top screenshot shows a handwritten diagram of a single neuron model. It has five input nodes labeled x_1, x_2, x_3, x_4, x_5 and a bias node b_0 . These inputs are connected to a summation node z . The output of z is passed through an activation function σ to produce the predicted output \hat{y} . The cost function J is shown to the right. The diagram is labeled "Forward model" and "Gradient Descent".

The bottom screenshot shows a table of data points, a small neuron diagram, and mathematical formulas. The table is as follows:

Expts	\vec{x}	y	\hat{y}	$J^{(i)}$
(1)	$x_1^{(1)}, x_2^{(1)}, x_3^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
(2)	$x_1^{(2)}, x_2^{(2)}, x_3^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
(m)	$x_1^{(m)}, x_2^{(m)}, x_3^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$

Next to the table is a small neuron diagram with input x and output \hat{y} . To the right, the cost function J is defined as:

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

$$J = - \sum_{i=1}^m J^{(i)}$$

The gradient of the cost function is also shown:

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w} \text{ - Chain Rule}$$

Welcome back. In the previous videos we saw our forward model of logistic regression. We also saw how we could use it in order to simulate the OR gate, NOR gate etcetera. And during all these cases what we were doing is we were sort of guessing for weights. Without actually doing the gradient descent, sort of by heuristic sort of rule of thumb guessing we figured out some weights for which we could replicate the OR gate. Now obviously that happened well because

we had only four data points and we just had to fit a classifying line between those four data points. In general, of course it is not really possible to just guess for a good line, which is why we look at how to calculate the gradient of the logistic regression which is will complete our loop.

So remember for logistic regression we are dealing with input nodes as usual, x_1, x_2, x_3 and we can have our usual bias unit sitting there. We have one output. Why only one output? Because this is binary classification, so the output is either going to be 1 or 0 or in our prediction case \hat{y} which lies between 0 and 1. And all these put together let us split it into two. We have our summation and then we have our sigmoid. This gives us \hat{y} . At the end of the summation, whatever we get is z and at the end of the sigmoid whatever we get is \hat{y} .

Now the weights which we do not know are here, w_0, w_1, w_2, w_3 and so on and so forth up till w_n . Remember just like in linear regression you could choose x_1, x_2 , to the x_1 square, x_3 to the x_1, x_2 , et cetera. So you can choose non-linear features as well. Okay, this is just as a reminder. So what is missing in this picture? We do not know these w s. This is just the forward model and we follow our usual procedure. You have the forward model. You guess for the w s. You get a \hat{y} , from the \hat{y} you get a J . From the J you feedback using $\frac{\partial J}{\partial w}$.

And through gradient decent or some version of gradient decent to calculate all that and improve guess for the w s. So which is missing here is of course this $\frac{\partial J}{\partial w}$. So I am going to write the whole thing in a vectorial representation. You have x vector, it runs through a sigma using a w , then through a non-linearity, through a sigmoid you get \hat{y} . And this \hat{y} is what gives you J . So just for simplicity or clarity, I am going to represent this slightly differently.

Let us take x vector, I will show the sigma here. I run it through w , I get z . I run z through a non-linearity and I get \hat{y} . From \hat{y} I get J . And what I want is $\frac{\partial J}{\partial w}$ vector. So question is, what is the expression for $\frac{\partial J}{\partial w}$ vector? Now this is fairly straightforward. Let us trace this back. $\frac{\partial J}{\partial w}$ vector equal to $\frac{\partial J}{\partial \hat{y}}$ times $\frac{\partial \hat{y}}{\partial z}$ times $\frac{\partial z}{\partial w}$ vector. So this is sort of the dependency. J changes because w changes. Why? Because J changes due to changes in \hat{y} . \hat{y} changes due to changes in z and z changes due to changes in w .

In other way, when I perturb w , it will perturb z which will perturb \hat{y} which will perturb J . So that is what we are chasing down here. This is essentially the Chain Rule. When we come to neural

networks, you will see that it is exactly the same idea which is applied for what is called a back propagation algorithm. So let us now calculate the expression for this. So let us look at each of these terms individually. First, what was J ? J if you recall, I had written as summation of $y \ln y$ hat, plus $(1 - y) \ln (1 - y)$ hat.

I am going to put an i here and i equal to 1 to m . And what is this i equal to 1 to m ? It is very very similar to what we did, in fact it is the same to what we did when we were doing linear regression. Recall, suppose you have x vector, as in the example that we have shown here, and you have multiple examples or multiple data points. So for the first data point x vector was x_1, x_2, \dots, x_n where you have n features. And I am going to put a superscript 1 to say that this is the first set of data points.

In the case of the OR example, we had four such data points. So you had an x_{11}, x_{21} which was 00 in that case. So similarly have $x_{12}, x_{22}, \dots, x_{n2}$. And we could have m such examples very similar to the multiple linear regression example that we did. Now we have the ground truth and since this is binary classification, ground truth here is y_1, y_2, \dots, y_m . y_1 will be either 0 or 1, y_2 will be either 0 or 1 and you have y_m . Then we have our y hat which is our prediction which is h of x given the parameters w .

So you are going to have $y_{hat 1}, y_{hat 2}$ up till $y_{hat m}$. Finally I am going to introduce something new just for clarity. I can have J_i , what is that? This is the loss just due to the i th example. So going back here, suppose you had four points and let us say this value was 1. This was y and y hat was let us say 0.8. The very fact that y hat and y differed will give you some amount of loss. So J_i is just the loss in the i th example. So you will have J_1, J_2, \dots, J_m .

Even in linear regression you can see if you have a line and you have multiple points, the loss for each one of those points will be denoted as J_i . So coming back here, all I am writing is J is the summation of all the binary entropy losses from each individual data point.

(Refer Slide Time: 9:31)

The image shows a blackboard with the following handwritten equations:

$$J = \sum_{i=1}^n J^{(i)}$$

$$\frac{\partial J}{\partial w} = \sum_{i=1}^n \frac{\partial J^{(i)}}{\partial w}$$

$$\frac{\partial J^{(i)}}{\partial w} = \frac{\partial J^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z} \frac{\partial z}{\partial w} \quad \hat{y} = \sigma(z)$$

$$J = -\left\{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right\}$$

$$\frac{\partial J}{\partial \hat{y}} = -\left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\} \quad \frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z}$$

So it is simply the summation of all individual losses. So suppose I want $\frac{\partial J}{\partial w}$, this is also going to be, actually I will remove the minus, let us include the minus within the J_i . This is $\frac{\partial J_i}{\partial w}$ vector. So we have this $\frac{\partial J_i}{\partial w}$ vector is equal to $\frac{\partial J_i}{\partial \hat{y}}$, $\frac{\partial \hat{y}}{\partial z}$, $\frac{\partial z}{\partial w}$ vector.

Technically speaking, this is \hat{y} . Okay, for now I am going to drop the i in the future expressions and we will just sum it back just for clarity of notation. Let us look at this term first. So I am going to drop the i as I said. J is $-\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\}$. Notice that the derivative is with respect to \hat{y} because that is what depends on w . y is fixed, y is what we gave, y are the labels that we give. \hat{y} is the parameter that, \hat{y} is the hypothesis function that we get out of the prediction using our parameter w .

So what is $\frac{\partial J}{\partial \hat{y}}$? $\frac{\partial J}{\partial \hat{y}}$, this term, is $-\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\}$. So we have that expression. What about $\frac{\partial \hat{y}}{\partial z}$? $\frac{\partial \hat{y}}{\partial z}$ is, remember \hat{y} is simply sigmoid of z , that is how we calculated it.

(Refer Slide Time: 11:46)

Grand Total Prediction = $f(x; w)$

Expts	\vec{x}	y	\hat{y}	$J^{(i)}$
(1)	$x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
(2)	$x_1^{(2)}, x_2^{(2)}, \dots, x_n^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
(m)	$x_1^{(m)}, x_2^{(m)}, \dots, x_n^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$

Diagram: $x \rightarrow z \rightarrow \hat{y} \rightarrow J$

$\frac{\partial J}{\partial \hat{y}} = ?$

$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}$ - Chain Rule

$J = -\sum_{i=1}^m \left[y^{(i)} \ln y^{(i)} + (1-y^{(i)}) \ln (1-y^{(i)}) \right]$

$J = \sum_{i=1}^m J^{(i)}$

$\frac{\partial J}{\partial w} = \sum_{i=1}^m \frac{\partial J^{(i)}}{\partial w}$

$\hat{y} = 0.8$

If you come here to this picture, \hat{y} was calculated as g of z or sigmoid of z .

(Refer Slide Time: 11:55)

$\frac{\partial J^{(i)}}{\partial \hat{y}} = \frac{\partial J^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z} \frac{\partial z}{\partial \hat{y}}$ $\hat{y} = \sigma(z)$

$J = -\left[y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right]$

$\frac{\partial J}{\partial \hat{y}} = -\left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\}$ $\frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma'(z) = \sigma(z) (1-\sigma(z)) = \hat{y} (1-\hat{y})$

$\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = -\left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\} \hat{y} (1-\hat{y})$

Now luckily in a previous video we already calculated this derivative. This was simply sigmoid of z times 1 minus- sigmoid of z , which is the same as \hat{y} times 1 minus- \hat{y} . So if you put it together, then you get $\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z}$ multiplied by $\frac{\partial \hat{y}}{\partial z}$, is equal to minus- $\frac{y}{\hat{y}}$ minus- $\frac{1-y}{1-\hat{y}}$ times \hat{y} times 1 minus- \hat{y} .

(Refer Slide Time: 12:37)

$$J = -\left\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\right\}^2$$

$$\frac{\partial J}{\partial \hat{y}} = -\left\{\frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}\right\} \quad \frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z)) = \hat{y}(1-\hat{y})$$

$$\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = -\left\{\frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}\right\} \hat{y}(1-\hat{y})$$

$$\frac{\partial J}{\partial z} = -\underbrace{\{y - \hat{y}\}}_{\text{-Error}} \quad (\text{Check!})$$

And if you calculate this, you get a very simple expression which is simply minus- y minus- y hat. Please check this for yourself. So we can call this expression even del J del z by the chain rule. Del J del z is equal to minus- y minus- y hat. This is simply the negative of the error. That is if you made a prediction y, y hat and the ground truth was y, del J del z is simply minus- y minus- y hat.

(Refer Slide Time: 13:16)

$$\frac{\partial J}{\partial w} = \sum_{i=1}^n \frac{\partial J^i}{\partial w}$$

$$\frac{\partial J^i}{\partial w} = \frac{\partial J^i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial z} \quad \hat{y} = \sigma(z)$$

$$z = \vec{w}^T \cdot \vec{x} = [1 \ x_1 \ \dots \ x_n]$$

$$\text{includes } w_0 = w^T x = x^T w.$$

$$J = -\left\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\right\}^2$$

$$\frac{\partial J}{\partial \hat{y}} = -\left\{\frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}\right\} \quad \frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z)) = \hat{y}(1-\hat{y})$$

$$\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = -\left\{\frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})}\right\} \hat{y}(1-\hat{y})$$

We have one other term left in this expression which is $\frac{\partial z}{\partial w}$. Remember z was vector dotted with x where w vector includes w not. This expression we saw the augmented w and the augmented x . x now includes 1, x_1 , up till x_n . So this can be written as $w^T x$, also as $x^T w$.

(Refer Slide Time: 14:08)

The image shows a blackboard with the following handwritten equations:

$$\frac{\partial J}{\partial \hat{y}} = - \left\{ \frac{y}{\hat{y}} - \frac{(1-\hat{y})}{(1-\hat{y})} \right\} \hat{y}(1-\hat{y})$$

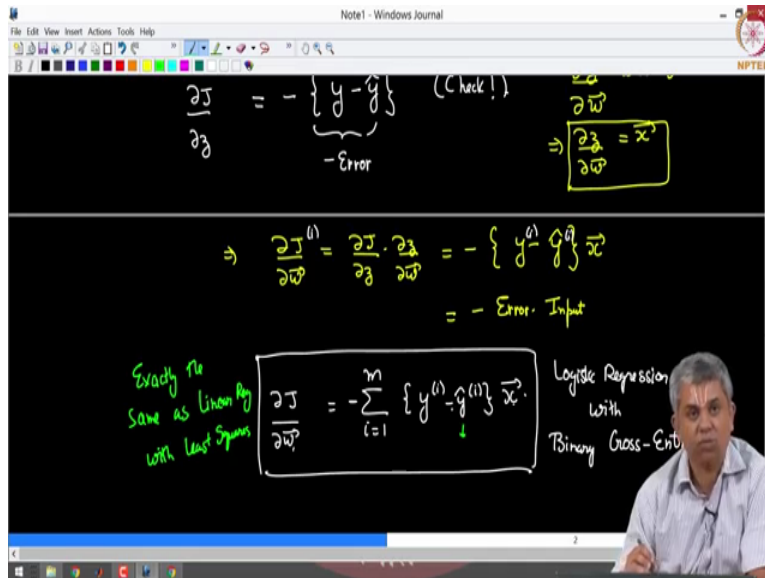
$$\frac{\partial J}{\partial \hat{y}} = - \underbrace{\{y - \hat{y}\}}_{\text{Error}} \quad (\text{Check!})$$

On the right side, there is a note: $\frac{\partial z}{\partial w}$ where $z = x^T w$. Below this, it is boxed: $\Rightarrow \frac{\partial z}{\partial w} = x$.

At the bottom, the final result is derived: $\Rightarrow \frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w} = - \{y - \hat{y}\} x = - \text{Error} \cdot \text{Input}$

We saw in the linear algebra videos that in such a case when you have $\frac{\partial J}{\partial w}$ where z is equal to $x^T w$, this means $\frac{\partial z}{\partial w}$ is simply equal to x vector. You can do this by algebra also or you could do this using the notation that I used in week 1 of this course. So if you put these two together, you get $\frac{\partial J}{\partial w}$, equal to $\frac{\partial J}{\partial z}$, multiplied by $\frac{\partial z}{\partial w}$. This is equal to minus- y minus- \hat{y} , multiplying x vector, which is minus- the error multiplied by the input. Now there are several noteworthy things here. I will go over them one by one.

(Refer Slide Time: 15:13)



First I will the general expression, this of course is $\frac{\partial J}{\partial w}$. So if I look at $\frac{\partial J}{\partial w}$, this is going to be summation from i equal to 1 to m of y_i , minus- \hat{y}_i , multiplied by x vector. So for example, if you wanted $\frac{\partial J}{\partial w}$ not, you will get y_i minus- \hat{y}_i , multiplied by x not which is simply 1, so on and so forth. For each component of w , you take the corresponding component of x . This is a vector, that is also a vector. So this expression is for logistic regression with binary cross entropy. Now some of you might recall that we had exactly the same expression, except for the factor of 1 by m which we can either include or not, exactly the same as linear regression with least squares.

So please refer back to your notes and check whether this is true or not. So this is actually remarkable. You get the same expression for the gradient for logistic regression with binary cross entropy as you get with linear regression with least squares. This is true even when you start including regularization. So when you include the regularization of, as Dr Ganpathi had shown earlier, if you include regularization of λ norm of w square, that would also add normally within logistic regression also.

The same thing holds true even for neural networks. So but the expression is actually exactly the same. This does not mean that the J is the same. Remember that this \hat{y} means different things for linear regression and for logistic regression. For linear regression \hat{y} was simply w

transpose x , for logistic regression y hat is now sigmoid of w transpose x . So please do remember that.

(Refer Slide Time: 17:48)

= - Error Input

Exactly the same as Linear Reg with Least Squares

$$\frac{\partial J}{\partial w} = - \sum_{i=1}^m \{ y^{(i)} - g^{(i)} \} \vec{x}^i$$

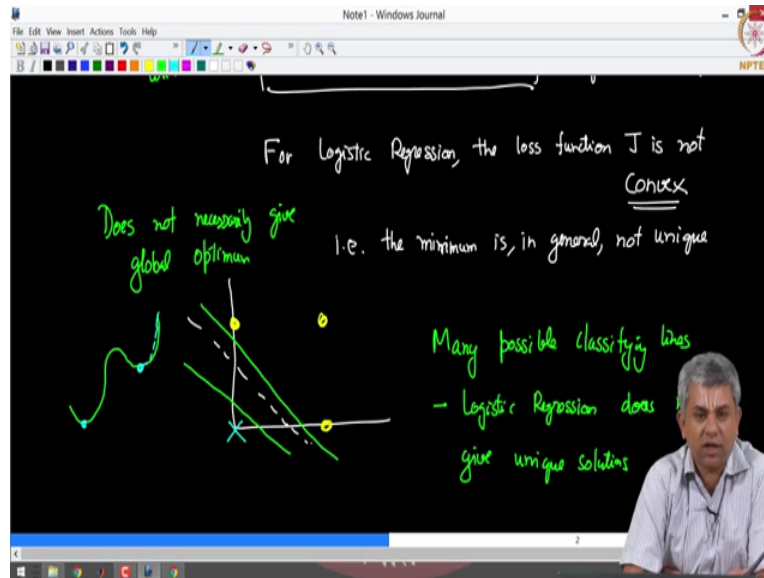
Logistic Regression with Binary Cross-Entropy

For Logistic Regression, the loss function J is not convex

i.e. the minimum is, in general, not unique

The second noteworthy thing here which is not obvious is that for logistic regression the loss function, J is not convex. So I will not go into what a convex set means but you can think when we had least squares generally the cost function will look simply like a paraboloid. So this is what is called a convex function where you have only 1 minima. The most important thing for us to know is that the minima, the minimum is in general not unique. We can understand this kind of intuitively too.

(Refer Slide Time: 18:48)



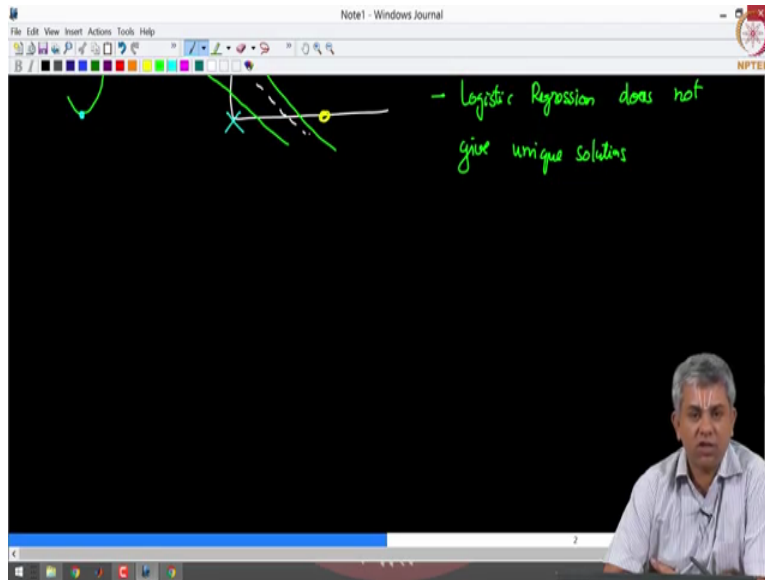
For example, we were trying to classify the set. Now notice when we were doing OR gate, we chose one line which was like this. But there is no reason to choose only this line. Even this line functions as a classifying line, this line functions as a classifying line. So there are many possible classifying lines. So logistic regression does not give unique solutions. So in our whole process depending on which w you start with, remember that our initial sets of weights we were guessing randomly. So the random guess depending on where it is, you might get one classifying line or the other.

So logistic regression actually depends on your initial conditions which classifying line you get. So that is an important constraint. This is also a problem with neural networks. It will never necessarily, does not necessarily give the global optimum. So in gradient decent let us say you have two local minima. At one place your gradient is 0, at another place also your gradient is 0. For example, if you have something of this sort, it is quite possible that your gradient decent comes here and gets stuck and does not move from there rather than come here which would actually be the local optimum.

So this is true whether it is logistic regression or whether it is neural networks which we will see shortly in future videos. But for linear regression typically if you converge, you will converge to only the global optimum because there is only one optimum because it is a convex function. So

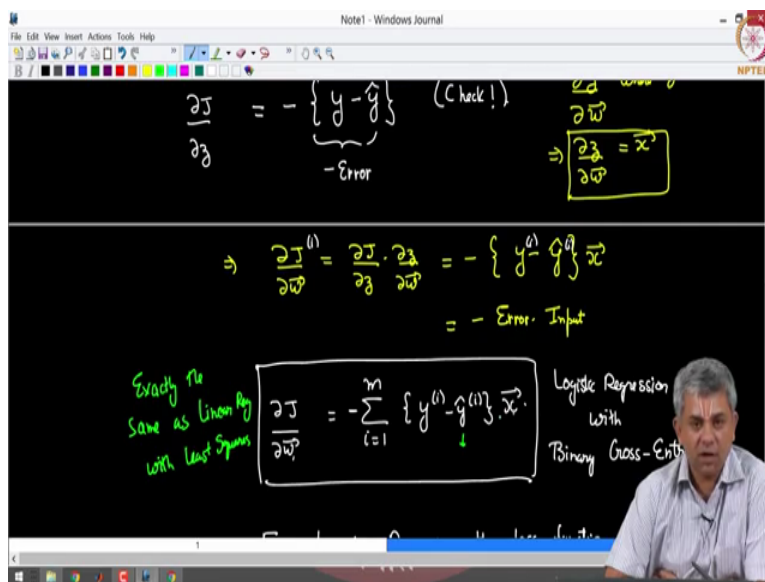
this is important for you to remember that logistic regression does not necessarily deal with convex loss functions.

(Refer Slide Time: 21:15)



So in order to sort of solidify your understanding of the logistic regression process, let us now look at a code which does logistic regression. You will see that this is just a small modification of the linear regression code that we had already used earlier. We had written a generalized linear regression code.

(Refer Slide Time: 21:40)



So since the gradient expression for logistic regression is practically indistinguishable from that of linear regression we will use the same code in order to do logistic regression. I will take the simple example of the OR gate.

(Refer Slide Time: 21:50)

- Logistic Regression does not give unique solutions

x	y
(0,0)	0
(0,1)	1
(1,0)	1
(1,1)	1

x_1, x_2

$$z = w \cdot x$$

$$\hat{y} = g(z)$$
 } vectorized expression

And in this case we will start with once again our example, 1, 2, 3, 4, we will treat it as four data points. And x vector which was 00, 01, 10, 11 and we have our ground truth which was 0, 1, 1, 1 and we put the same model as before, a summation followed by a sigmoid gives us \hat{y} after taking in x . Now what I will also be doing is both in code which you will notice as well as in the expression here, we will write it as if we are writing vector expressions just to tell you how that works.

So if you have x_1 , x_2 , and x_{not} here which is just 1, we sum these three up, w_{not} , w_1 , w_2 . What comes out is \hat{y} after you put sigmoid and g . You can write this as z which is the intermediate output of this, z is equal to $w \cdot x$. Remember this w also includes the w_{not} which is the biased term and x also includes x_{not} . And \hat{y} is equal to g of z . This would be the vectorized expression. And this is what lets us write a general code with great amount of ease because suppose I decide to use other non-linear features.

For example, you saw the XOR gate example, in the XOR gate example suppose you want to include x_1 square, x_2 square et cetera, you could simply include extra features x_3 , x_4 and that would logistic regression would still work as it is.

(Refer Slide Time: 24:08)

The image shows a blackboard with handwritten notes and diagrams. On the left, there is a truth table for an OR gate:

s	x	OR y
1	(0,0)	0
2	(0,1)	1
3	(1,0)	1
4	(1,1)	1

Below the table, the inputs are labeled x_1, x_2 . To the right of the table, there is a diagram of a single neuron with an input x , a summation node Σ , and an activation function g leading to output y . Below this, there is a diagram of a vectorized neuron with inputs x_1 and x_2 , weights w_1 and w_2 , a bias w_0 , a summation node Σ , and an activation function g leading to output y . The equations $z = \vec{w} \cdot \vec{x}$ and $y = g(z)$ are written below the vectorized neuron diagram, with a note "vectorized dot" next to the first equation.

So we will look at a vectorized code for this in the next video which would be the code. Thank you.