

Logistic Regression Gradient Descent

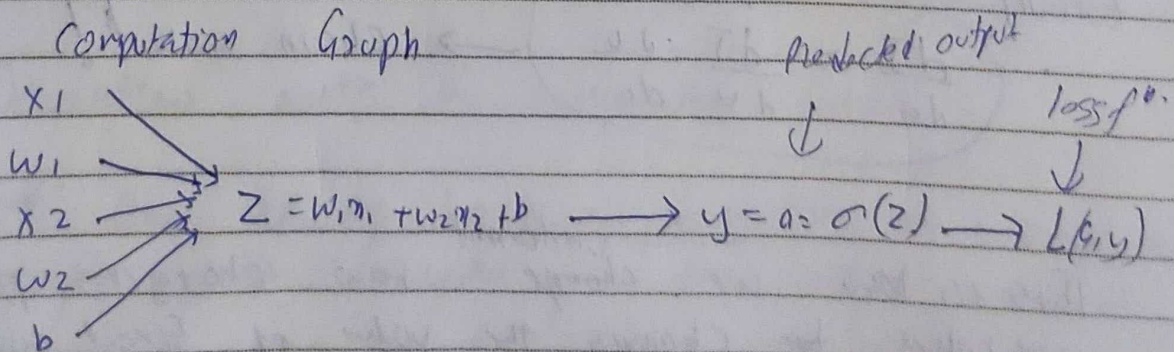
we're trying to study Logistic Regression
thru computation graphs

In LR.

$$z = w^T x + b$$

prediction $\hat{y} = a = \sigma(z)$

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$



now go backward to compute derivative of Loss
wrt to initial variables

basically how much J changes on changing a .

remembers own loss!

1st step

we enter $da = d \frac{L(a,y)}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$

2nd step (in going backward)

$$dz = \frac{dL(a,y)}{dz} = \frac{dL(a,y)}{da} \cdot \frac{da}{dz} \quad (\text{chain rule!!})$$

← remembers our a fn.

$$dz = \frac{a-y}{a(1-a)}$$

3rd step in going back.

$$\frac{\partial L}{\partial w_1} = dw_1 = \eta_1 dz \quad ; \quad dw_2 = \eta dz$$

\downarrow
 $\eta_1 (a-y)$

$$; \quad db = dz$$

\downarrow
 $(a-y) \quad \eta_2 (a-y)$

then we apply Gradient descent as

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

→ All of this was only for 1 training example,

Now we'll see how we apply Gradient Descent for multiple training examples.

Gradient Descent on m examples

$$\text{Cost } f^n := J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

we notice

so the overall cost function is an average of individual loss f^n from each training example

So

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial L(a^{(i)}, y^{(i)})}{\partial w_1}}_{\downarrow}$$

we calculate this part

So now we just have to average them !!

So :- for model with 2 input parameters and 1 middle layer
start with $J=0$; $dw_1=0$; $dw_2=0$; $db=0$

For $i=1$ to m :

1 step
of
Gradient
descent

Forward propagation

$$\begin{cases} z^{(1)} = w^T x^{(1)} + b \\ a^{(1)} = \sigma(z^{(1)}) \\ J += -[y^{(1)} \log(a^{(1)}) + (1-y^{(1)}) \log(1-a^{(1)})] \end{cases}$$

Backward propagation

$$\begin{cases} dz^{(1)} = a^{(1)} - y^{(1)} \\ dw_1 += x_0^{(1)} dz^{(1)} \\ dw_2 += x_1^{(1)} dz^{(1)} \\ db += dz^{(1)} \end{cases} \quad \left. \begin{array}{l} \text{up only two} \\ n=2. \end{array} \right\}$$

outside of for loop

$$J /= m; \quad dw_1 /= m; \quad dw_2 /= m; \quad db /= m$$

now dw_1, dw_2, db represent the respective changes we need to make to reduce cost function

$$\begin{aligned} \text{So we do } w_1 &:= w_1 - \alpha dw_1 \\ w_2 &:= w_2 - \alpha dw_2 \\ b &:= b - \alpha db \end{aligned}$$

This completes one step of Gradient descent
* Now we repeat this process until we find Principal minima using Gradient Descent

Now, Guess what — this is very bad.
this is computationally stupid.

we're ~~not~~ ^{not}

In 1 step of Gradient descent we're running $O(mn)$ Time Complexity and then we repeat Gradient descent n times to reach minima
 \downarrow
 m features \times updating n parameters

running nested for loops in Deep Learning algo's is considered a bad practice

So we need to find a different way
 \rightarrow Vectorization