

Derivatives #

The derivative formula predicts that if you nudge a by a tiny little bit then

$y = f(a) = a^3$ should go up

by $\frac{df(a)}{da} \cdot 3a^2$ and indeed it does -
da

Ex 2: if you bump up a by 0.001 you can expect

$f(a)$ to go up by $\frac{1}{2}$ as much as a

so if $f(a) = \log$

the increase in y is 3 times the increase in x

Final output variable that I really care about

$$\frac{dy}{dx} = da \text{ in code}$$

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

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

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

→ My job when I implement logistic regression, is to learn parameters, w and b so that \hat{y} becomes a good estimate of the chance of y being equal to one.

$$\text{given } x \text{ want } \hat{y} = p(y=1|x)$$

If $y=1$, we wanna make \hat{y} large
 $\hat{y} \approx 1$

If $y=0$, we wanna make \hat{y} small
 $\hat{y} \approx 0$

Binary classification:
output either 0 or 1

$$\underline{\mathcal{L}} = "d\omega_1" = \omega_1 \cdot dz$$

$$\rightarrow \omega_1 \quad db = dz$$

$$d\omega_2 = \omega_2 \cdot dz$$

$$\underline{\mathcal{L}} = \omega_1 z_1 + \omega_2 z_2 + b$$

$$[a = \sigma(z)] = \mathcal{L}(a, y)$$

$$w_1 := \omega_1 - \alpha d\omega_1$$

$$w_2 := \omega_2 - \alpha d\omega_2$$

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

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

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

Algorithm

$$J=0; \text{ dev}_1=0; \text{ dev}_2=0; db=0$$

for $i=1$ to m {

$$Z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(Z^{(i)})$$

$$J+ = [y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dZ^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1+ = x_1^{(i)} dZ^{(i)} \quad | \quad n=2$$

$$dw_2+ = x_2^{(i)} dZ^{(i)} \quad | \quad n=2$$

$$db+ = dZ^{(i)} \gamma$$

$J \neq m$

$$dw_1/m; dw_2/m; db/m$$

compute w^T

$$\boxed{Z = w^T x + b}$$

$$Z=0$$

for i in range($\hat{\Omega}_w$): {

$$z \doteq w[i] \otimes x[i]^T$$

$$Z \doteq ls$$

Numpy $z = np.\text{dot}(w, x) + ls$

a = $np.\text{array}([1, 2, 3, 4])$

perim(a) $\left[\begin{matrix} \text{Shift} \\ +^2 \text{ and } -j \end{matrix} \right]$

import time

a = $np.\text{random.rand}(1000000)$

million dimensional array

b = $np.\text{random.rand}(1000000)$

tic = time.time()

c = $np.\text{dot}(a, b)$

toc = time.time()

v = $np.\text{zeros}((n, 1))$

v[i] = math.exp(v[i])

↳ trainings inputs stacked together
in different columns like this

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}$$

$$w^T X + [b \ b \ \dots \ b]$$

$| \times m$

$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b \\ \dots & w^T x^{(m)} + b \end{bmatrix}$$

$$Z = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix}$$

$$Z = np.\text{dot}(w.T, X) + b$$

$$A = \underline{\begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix}}$$

$$dz^{(1)} = a^{(1)} - y^{(1)}$$

$$dz^{(2)} = a^{(2)} - y^{(2)} \dots$$

$$dz^{(m)} = a^{(m)} - y^{(m)}$$

$$dZ = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix}$$

$$dZ = A - Y$$

$$y = [y^{(1)} \dots y^{(m)}]$$

$$dw = 0$$

$$dw + = x^{(1)} dz^{(1)}$$

$$dw + = x^{(1)} dz^{(2)} \dots \dots \quad dw + = x^{(m)} dz^{(m)}$$

$$db = 0$$

$$db + = dz^{(1)}$$

$$db + = dz^{(2)} \dots -$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \text{ np.sum } (dz)$$

$$dw = \frac{1}{m} X^T dz^T$$

$$= \frac{1}{m} \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \begin{bmatrix} x_1^{(1)} dz^{(1)} & \dots & x_1^{(m)} dz^{(m)} \end{bmatrix}$$

get rid of Biggar for loop

$$Z = w^T X + b$$
$$= np \cdot \text{dot}(w.T, X) + b$$

$$A = \mathcal{T}(Z)$$

$$dZ = A - Y$$

$$\underbrace{dw = \frac{1}{m} X dZ^T}_{\text{gradient}}$$

$$dw = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} np \cdot \text{sum}(dZ)$$

gradient descent update

$$w := w - \alpha dw$$

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

$$cal = A * \text{sum}(\text{axis}=0)$$

$$\text{fmin}(cal)$$

$$\text{percentage} = 100 * A / cal \cdot \text{reshape}(1, 4)$$

$$a = np \cdot \text{random.random}((5, 1))$$

$$\text{array}(a \cdot \text{shape} = (5, 1))$$

$a = a \cdot \text{reshape}((5, 1))$

$$\overline{p(y|x)} = y^2 (1-y)^{(1-y)}$$

IID = Identity independently distributed

P(labels in training set)

$$= \prod_{i=1}^m p(y^{(i)} | x^{(i)})$$

$$\log P(\dots) = \sum_{i=1}^m \log p(y^{(i)} | x^{(i)})$$

Quiz - 2

(10) computation graph

$$J = u + v - w$$

$$= a * b + a * c - (b + c)$$

$$= a(b + c) - (b + c)$$

$$= (a - 1)(b + c)$$

gradient = vector quantity
derivative = scalar quantity

gradient of a sigmoid function

sigmoid_derivative(x)
 $= \sigma'(x) = \sigma(x)(1 - \sigma(x))$

$$\frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right)$$

grad \cong slope \cong derivative

grad descent converges faster
after normalization

$x_norm = np.linalg.norm$

(x , ord=2, axis=1, keepdims
=True)

$x_sum = np.sum(x_exp,$
axis=1, keepdims=True)

dot product of a row
vector & row vector is [1]

$$\textcircled{9} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} * \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$\textcircled{8} \quad \begin{bmatrix} 2 & 1 & 1 & 1 \\ 1 & 2 & 1 & 1 \\ 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 2 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

O/p of a neuron is
 $a = g(Wx + b)$ where g is
 the activation function

(sigmoid, tanh, ReLU)

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

tanour as logistic function

math lib when inputs of functions are real numbers

But when inputs of func are matrices & vectors
 numpy lib more useful

The loss is used to evaluate performance of model.

Bigger loss = more difference predictions (\hat{y}) from true values (y)

Gradient Descent = optimization algorithm = train model & minimize cost

if $x = [x_1, x_2, \dots, x_n]$
then np. dot (x, x)
 $= \sum_{j=0}^n x_j^2$

X_.flatten = X_.reshape
(X_.shape[0], -1).T
(a, b, c, d) \rightarrow (b * c * d, a)

Standardize the dataset

3 Subtract the mean of the whole numpy array from each example, and then divide each example by the standard deviation of the whole numpy array

Common steps for preprocessing a new dataset

- (i) figure out the dimensions and shapes of the problem (m_{train} , m_{test} , num_px , ...)
- (ii) reshape the datasets such that each example is now a vector of size $(\text{num_px} * \text{num_px}^3, 1)$
- (iii) “Standardize” the data

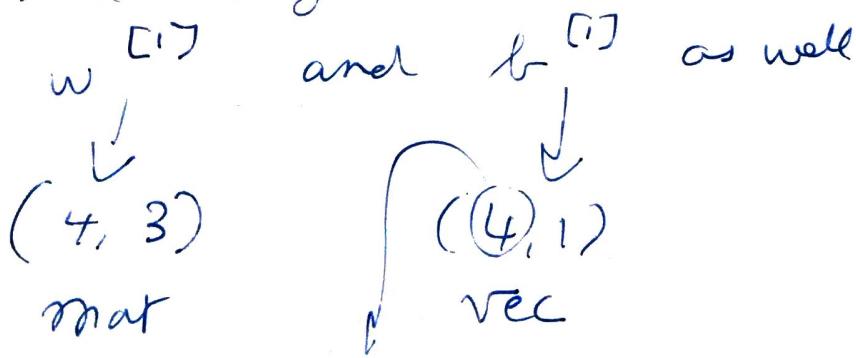
Learning rate α determines how rapidly we update the parameters

2 Overfitting: training accuracy much larger than test accuracy

Week 3

3 Hidden layer not visible values in the training set.

In a layer we got



4 hidden units in the layer and 3 input layer features

So the hidden layer has associated with it these parameters

Ques Stack the nodes vertically

To form the vector $\hat{z}^{(l)}$

$$\hat{z}^{(l)} = \begin{bmatrix} -w_1^{(l)T} \\ -w_2^{(l)T} \\ -w_3^{(l)T} \\ -w_4^{(l)T} \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix}$$

$$+ \begin{bmatrix} b_1^{(l)} \\ b_2^{(l)} \\ b_3^{(l)} \\ b_4^{(l)} \end{bmatrix} = \begin{bmatrix} w_1^{(l)T}x + b_1^{(l)} \\ w_2^{(l)T}x + b_2^{(l)} \\ w_3^{(l)T}x + b_3^{(l)} \\ w_4^{(l)T}x + b_4^{(l)} \end{bmatrix}$$

$\hat{z}^{(l)}$

$$\hat{z}^{(l)} = w^{(l)T}x + b^{(l)}_{(4,1)}$$

$$a^{(l)} = \sigma(\hat{z}^{(l)})$$

$$x = a^{(0)}$$

$$\hat{y} = a^{(L)}$$

alias for vector of input features x

$$W^{[2]} \quad b^{[2]} \\ (1, 4) \quad (1, 1) \\ z^{[2]} \quad (1, 1) = (1, 4)(4, 1) \\ + (1, 1)$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

Superscript round bracket i
refers to the i^{th} training
example

Superscript square bracket i
refers to this stack of
nodes called a layer.

$$f_{3i} \stackrel{i=1 \text{ to } m}{=} f_{[3](i)} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[1](i)} = f_{[2](i)}$$

$$f_{[2](i)} = W^{[2]} a^{[0]} + b^{[2]}$$