

Deep Learning Specialization

week 2

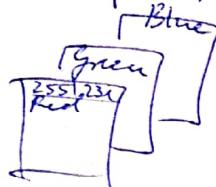
course I

Binary Classification

with
↳ Logistic Regression

[first image] → $\begin{cases} 1 (\text{cat}) \\ 0 (\text{non cat}) \end{cases}$ } output label y

How is image represented?



Pixels 64×64
64 image
64

union in feature vector x

$$x = \begin{pmatrix} 255 \\ 231 \\ \vdots \\ G \\ \vdots \\ B \end{pmatrix} \rightarrow \begin{array}{l} \text{total dimension} \\ 64 \times 64 \times 3 \\ L = 12288 \\ n_x = 12288 \end{array}$$

Notation

$$(x, y) \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

train, test, (McV...)

$$\bar{X} = \left[\begin{array}{c|c|c|c} \text{vectors} & & & \\ \hline X^{(1)} & | & X^{(2)} & | \\ \hline \vdots & | & \vdots & | \\ \hline X^{(m)} & | & \vdots & | \end{array} \right]_{n_x \times m}$$

Logistic Regression

Given X , we want $\hat{y}_{\text{act.}} = P(y=1|x)$
e.g. picture $x \in \mathbb{R}^{n_x}$ $0 \leq y \leq 1$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$ weights bias
w is n_x -dimensional vector b is real number

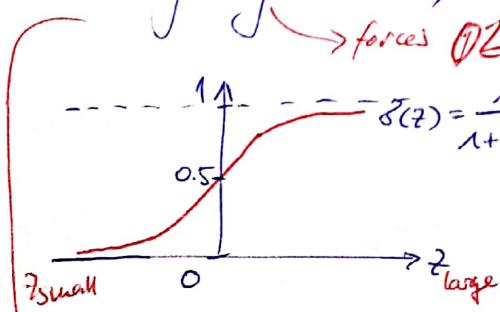
Output $\hat{y} = w^T x + b$ Linear Regression approach doesn't work because \hat{y} could be any number but has to be $0 \leq \hat{y} \leq 1$!

$$\hat{y} = g(w^T x + b)$$

\rightarrow forces $0 \leq \hat{y} \leq 1$!

$$\dots \quad \hat{y} = g(z) = \frac{1}{1+e^{-z}} \quad \text{If } z \text{ is large } g(z) \approx \frac{1}{1+0} = 1$$

$$\text{If } z \text{ is low } g(z) \approx \frac{1}{1+e^{-z}} \approx \frac{1}{1+b_0} \approx 0$$



\rightarrow first solve linear term, then activation!

rows columns

$X \in \mathbb{R}^{n_x \times m}$

$\hat{Y} = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$

$\hat{Y} \in \mathbb{R}^{1 \times m}$

Python command for finding shape of matrix

$X.\text{shape} = (n_x, m)$

$\hat{Y}.\text{shape} = (1, m)$

feature vector x and weight vector w
must have same dimension

not using this notation in this course (theta)

$X_0 = 1, X \in \mathbb{R}^{n_x \times 1}$

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

$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{pmatrix}_w$

parameters

keep separate!

$\hat{y}_{\text{raw}} = \hat{y}_0(x)$

$\hat{y} = g(\hat{y}_{\text{raw}})$

y

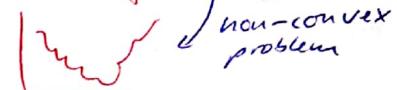
X

①

Logistic Regression cost function

$$\hat{y} = \sigma(w^T x + b) \text{ with } \sigma(z) = \frac{1}{1+e^{-z}}$$

Given $\{(x^{(i)}, y^{(i)}) - (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$ i -th example

Loss (error) function: $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2 \rightarrow \text{don't do in logistic regression}$ 

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

If $y=1$:= $L(\hat{y}, y) = -\log \hat{y} \leftarrow$ want as small as possible (as possible)

If $y=0$:= $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $\log(1-\hat{y})$ as large, want \hat{y} as small as possible

$$\text{Cost function: } J(w, b) = \frac{1}{m} \sum L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$

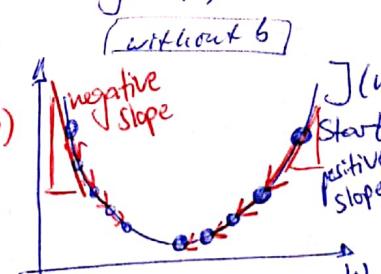
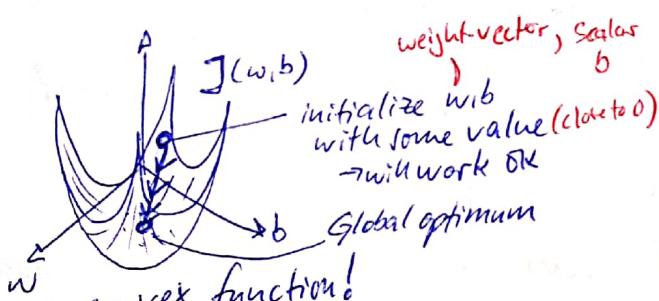
Loss function computes the error for a single training example, the cost function is the average of the loss functions of the entire training set

Gradient Descent (optimization)

Recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$

cost function averaging loss function \rightarrow Want to find w, b that minimize $J(w, b)$



Repeat {

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

$\underbrace{\frac{\partial J(w)}{\partial w}}_{\text{gradient term}}$

}

learning rate

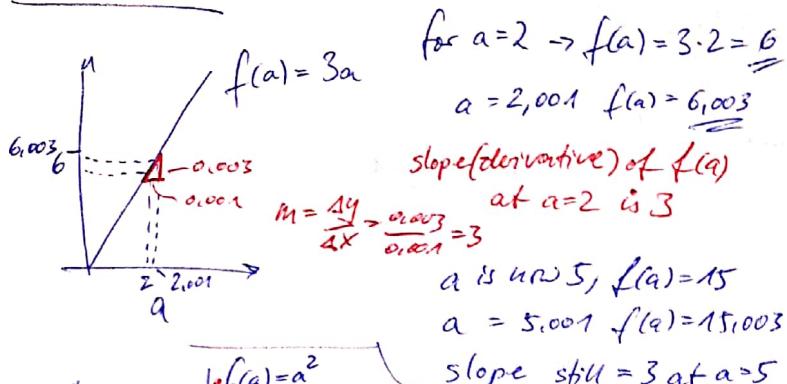
$$\frac{\partial J(w)}{\partial w} = ?$$

$$J(w, b) \quad w := w - \alpha \frac{\partial J(w, b)}{\partial w} \quad \text{partial derivatives}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b} \quad \text{in code "db"}$$

derivative is slope of a function

Derivatives

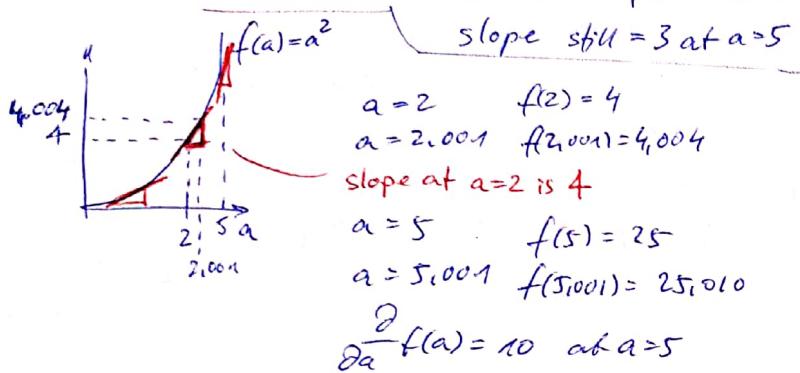


$$\frac{\partial f(a)}{\partial a} = 3 = \frac{\partial}{\partial a} f(a)$$

formal definition:

slope after infinitesimal shift?

→ slope always same for linear function



$$\frac{\partial}{\partial a} f(a) = \frac{\partial}{\partial a} a^2 = 2a$$

$$f(a) = a^2 \quad \frac{\partial}{\partial a} f(a) = 2a \quad a=2 \quad a=2,001 \\ \hookrightarrow f(a)=4 \quad \hookrightarrow f(a) \approx 4,004$$

$$f(a) = a^3 \quad \frac{\partial}{\partial a} f(a) = 3a^2 \quad a=2 \quad a=2,001 \\ \hookrightarrow f(a)=8 \quad \hookrightarrow f(a) \approx 8,012$$

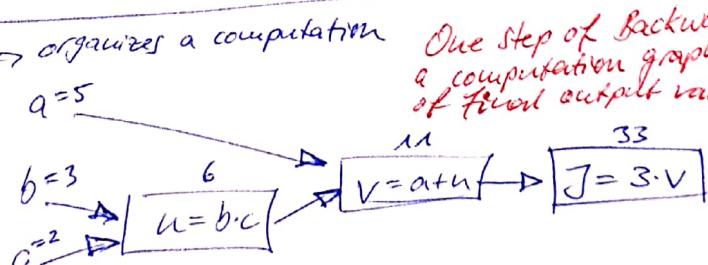
$$f(a) = \log(a) \quad (\ln(a)) \quad \frac{\partial}{\partial a} f(a) = \frac{1}{a} \quad a=2 \rightarrow f(a)=0,69315 \\ a=2,001 \rightarrow f(a)=0,69365 \quad \Delta=0,0005$$

$\frac{\partial}{\partial a} f(a) = \frac{1}{2}$

Computation Graph

$$J(a, b, c) = 3(a + bc)$$

- compute $u = bc$
- compute $v = a + u$
- compute $J = 3 \cdot v$



Derivatives with a Computation Graph

(Forward propagation)

$$a = 5$$

$$b = 3$$

$$c = 2$$

$$\frac{\partial J}{\partial a} = ? = 3$$

$$\frac{\partial J}{\partial a} = ? = 3$$

$$\begin{aligned} \frac{\partial J}{\partial v} &= ? = 3 \\ \frac{\partial v}{\partial a} &= 1 \\ \frac{\partial J}{\partial a} &= 3 \end{aligned}$$

Chain rule!

cost function

$$33$$

$$v = a + u$$

$$J = 3 \cdot v$$

J defined as $J = 3v$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\begin{aligned} a &= 5 \rightarrow 5.001 \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

) changes
↓ changes
↓ changes

Focal output variable, usually target for optimization
in this case is \underline{J}

In this course

$$\frac{\text{dFinal Output Var}}{\text{dVar}}$$

$\frac{\text{dVar}}{\text{dVar}}$
Derivative of final output variable J , i.e.
w.r.t. to various intermediate quantity

Backpropagation

$$a = 5$$

$$\frac{\partial J}{\partial a} = da = 3$$

$$\frac{\partial J}{\partial b} = db = 6$$

$$\frac{\partial J}{\partial c} = dc = 9$$

$$\frac{\partial J}{\partial v} \frac{\partial v}{\partial u} = -3$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial c} = 9$$

$$\begin{aligned} u &= b - c \\ v &= a + u \\ J &= 3v \end{aligned}$$

$$\begin{aligned} u &= 6 \rightarrow 6.001 \\ v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned}$$

$$b = 3 \rightarrow 3.001$$

$$u = 6 \rightarrow 6.002$$

$$v = 11 \rightarrow 11.002$$

$$J = 33.006$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \frac{\partial u}{\partial b} = 6$$

$$3 \cdot 2$$

Vectorization \rightarrow get rid of for-loops Whenever possible avoid explicit for-loops

$$z = \omega^T x + b \quad \omega = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

non-vectorized
normal

$$z = 0$$

for i in range(n_x):

$$z += \omega[i] \cdot x[i]$$

$$z += b$$

450 ms

vs

vectorized (Python)

$$z = \underbrace{\text{np.dot}}_{\text{computes } \omega^T x}(w, x) + b$$

for 1 million-dimensional vectors a, b

1.5 ms

GPU SIMD

CPU Single instruction
multiple data

enable Python Numpy
for better parallelism

more examples

$$u = Av \quad \text{non-vectorized}$$

$$u_i = \sum_j A_{ij} v_j$$

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

for i ...

for j ...

$$u[i] += A[i][j] \cdot v[j]$$

vectorized

$$u = \text{np.dot}(A, v)$$

$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$ want to apply exponential operation on every element of a matrix/vector

$$u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

\rightarrow non-vectorized

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

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

vectorized

import numpy as np

~~for i in range(n):~~

$$u = \text{np.exp}(v)$$

np.log(v) \rightarrow elementwise log

np.abs(v) \rightarrow absolute values

np.maximum(v, 0)

$$v * x2$$

$$\sqrt{v}$$

Numpy
built-in
functions

logistic regression derivatives
 \rightarrow problem with 2 derivatives

$$dw = \text{np.zeros}((n_x, 1)) \rightarrow \text{initialization}$$

operation

$$dw += x^{(i)} d z^{(i)} \rightarrow 2^{\text{nd}} \text{ for loop}$$

$$dw /= m \rightarrow \text{averaging}$$

for loop looping over training examples remain

(6)

Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(m)} = w^T x^{(m)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(m)} = \sigma(z^{(m)})$$

matrix $\bar{X} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \\ 1 & 1 & \dots & 1 \end{bmatrix}_{n_x \times n_m}$

compute $z^{(1)} \dots z^{(m)}$ in 1 line of code:

$$\bar{Z} = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T \bar{X} + [b \ b \ \dots \ b]$$

$1 \times m$ matrix

row vector

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

$$= [w^T x^{(1)} + b, w^T x^{(2)} + b, \dots, w^T x^{(m)} + b]_{1 \times m}$$

one line of code that calculates \bar{Z} that contains all of lowercase z 's

$1 \times 1 \rightarrow R$

real number

↳ Python automatically expands the real number out to a $1 \times m$ row vector

↳ "Broadcasting"

Next want to find a way to compute

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(\bar{Z})$$

new variable
activations!

computes all lower case as

Vectorizing Logistic Regression's Gradient Computation

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

$$d\bar{Z} = [dz^{(1)} \ dz^{(2)} \ \dots \ dz^{(m)}] \quad \text{row vector}$$

$1 \times m$ matrix

$$A = [a^{(1)} \ \dots \ a^{(m)}] \quad Y = [y^{(1)} \ \dots \ y^{(m)}]$$

$$d\bar{Z} = A - Y = [a^{(1)} - g^{(1)} \ \dots \ a^{(m)} - g^{(m)}]$$

got rid of 1st for loop but still had

$$dw = 0$$

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

$$dw += x^{(2)} dz^{(2)}$$

$$\vdots$$

$$dw / m$$

$$db = 0$$

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

$$db += dz^{(2)}$$

$$\vdots$$

$$db / m$$

for multiple iterations still need for-loop

single iteration of gradient descent for log. Reg.

vectorized implementation

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

$$\text{Python: } \frac{1}{m} \cdot \text{np.sum}(dz)$$

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

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

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]$$

$n \times 1$ vector

$$Z = w^T X + b$$

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

$$A = \sigma(Z)$$

$$d\bar{Z} = A - Y$$

$$dw = \frac{1}{m} X d\bar{Z}^T \quad db = \frac{1}{m} \text{np.sum}(dz)$$

$$w := w - \alpha dw \quad b := b - \alpha db$$

FP and BP without using a for loop

(7)

Repeat :

for iter in range(1000): ← Still need for-loop for multiple iterations

$$\begin{aligned} Z &= w^T X + b \\ &= \text{np.dot}(w.T, x) + b \\ \text{implementation} \\ \text{of single} \\ \text{iteration} \\ \text{of gradient} \quad dZ &= A - Y \\ \text{Descent for} \quad dw &= \frac{1}{m} X dZ^T \\ \text{logistic} \quad db &= \frac{1}{m} \text{np.sum}(dZ) \\ \text{Regression} \quad w &:= w - \alpha dw \\ b &:= b - \alpha db \end{aligned}$$

} forward propagation
and back propagation
without using
a for-loop

Broadcasting in Python → technique to run code faster

Calories in 100g:

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

59 cal

Calculate % of calories from
carbs, protein, fat

Apple: $\frac{56}{59} \approx 94.9\%$

→ Can this be done without for-loop?

↪ $[\cdot \cdot \cdot] = A$
(3,4)

$A = \text{np.array}([[56, 0, 4.4, 68],$

$[1.2, 104, 52, 8],$

$[1.8, 135, 99, 0.9]])$

$cal = A \cdot \text{sum}(\text{axis}=0)$ want Python to sum vertically

apple beef eggs points

cal = [59, 229, 155, 76]

percentage = $A / \text{cal} \cdot \text{reshape}(1, 4) \times 100$ ← Python broadcasting

↪ [[99.9, 0, 2.83, 88.4]]

3x4 matrix divided

by 1x4 matrix

(cal is actually already
in right shape here)

↪ $A = \frac{3 \times 4 \text{ matrix}}{\text{cal } 1 \times 4 \text{ matrix}}$

more examples ↴

1
2
3
4

1
2
3
4

1
2
3
4

1
2
3
4

changes/expands
number into
vector of desired
dimension

↪ works for
row and line
vectors!!!

5) $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$

(1,n) → copy matrix
to (m,n)

c) $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 204 & 205 & 206 \end{bmatrix}$

⑧

Matlab/octave
→ lsxfun

Broadcasting in Python

General Principle / newer form of broadcasting

$$\begin{matrix} (m,n) \\ \text{matrix} \end{matrix} + \begin{matrix} (1,n) \\ \vdash \end{matrix} \rightarrow (m,n)$$

$$(m,1) + R$$

$$\left[\begin{matrix} 1 \\ 2 \\ 3 \end{matrix} \right] + 100 = \left[\begin{matrix} 101 \\ 102 \\ 103 \end{matrix} \right]$$

$$[1, 2, 3] + 100 = [101, 102, 103]$$

Note on Numpy vectors (bug avoidance)

expressibility

→ strength: flexibility

→ weakness: bugs, subtle bugs

You'll use "rank-1 vectors" arrays

example a):

$$a = np.random.rand(5) \quad \text{rather do}$$

$$\begin{matrix} \rightarrow \\ \text{print}(a.shape) \rightarrow (5,) \\ \text{print}(a.T) \rightarrow ... \end{matrix} \quad \begin{matrix} \text{rather not use} \\ \text{"rank 1 - array"} \\ \text{!} \end{matrix} \quad \begin{matrix} \rightarrow \\ \text{a.shape} \rightarrow (5,1) \end{matrix}$$

If not sure of dimension of one of the vectors → assert(a.shape == (5,1))
Inexpensive to execute

$a = np.random.rand(5,1) \rightarrow$ column vector
or (1,5) would create a row vector

$a = np.random.rand(5,1) \rightarrow$ column vector
 $a = np.random.rand(1,5) \rightarrow$ row vector

If you have a "rank-1 array" do → $a = a.reshape((5,1))$
(or (1,5))

(optional) Explanation of logistic regression cost function

Cost on m examples

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

$$\log p(\dots) = \sum_{i=1}^m \underbrace{\log(y^{(i)}|x^{(i)})}_{-L(y^{(i)}, y^{(i)})}$$

In statistics: principle of maximum likelihood estimation

choose parameters that maximize

$$-\sum_{i=1}^m L(y^{(i)}, y^{(i)})$$

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

average so results are better to scale

$$\text{summarize: } P(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

$$\text{If } y=1 \rightarrow \hat{y} \cdot \underbrace{(1-\hat{y})^{1-1}}_{=1} = \hat{y}$$

$$\text{If } y=0 \rightarrow 1 \cdot \underbrace{(1-\hat{y})^{1-0}}_{1-\hat{y}} = 1-\hat{y}$$

$$\begin{aligned} \log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{1-y} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -L(\hat{y}, y) \end{aligned}$$

→ Minimizing L corresponds to maximizing $\log p(y|x)$

test question

$$\begin{array}{l}
 \begin{array}{c}
 \cancel{a} \rightarrow u = a \cdot b \\
 \cancel{b} \rightarrow v = a \cdot c \\
 \cancel{c} \rightarrow w = b + c
 \end{array}
 \Rightarrow J = u + v - w \\
 = a \cdot b + a \cdot c - (b + c) \\
 = a \cdot b + a \cdot c - b - c
 \end{array}$$

$$J = (c-1) \cdot (b+a) = b \cdot c + a \cdot c - b - a$$

$$J = (a-1) \cdot (b+c) = a \cdot b + a \cdot c - b - c$$

$$x = [x_1, x_2, \dots, x_n]$$

$$\text{up. dot}(x, x) = \sum_{j=0}^m x_j^2$$

$$L_2(y, \hat{y}) = \sum_{i=0}^m (y^{(i)} - \hat{y}^{(i)})^2$$

$$\hookrightarrow \text{up.sum}(\text{up.dot}((y - \hat{y}).\text{abs}), (y - \hat{y}).\text{abs}))$$

don't really need as up dot
summing already

$$L_1(y, \hat{y}) = \sum_{i=0}^m |y^{(i)} - \hat{y}^{(i)}|$$

$$\hookrightarrow \text{up.sum}(\text{abs}(y - \hat{y}))$$

image \rightarrow 3D array of shape (length, height, depth = 3)

\hookrightarrow convert to vector of shape (length · height · 3, 1)

it wants: array v of shape (a, b, c) in to shape (a · b · c, 1)

$$\hookrightarrow v = v.\text{reshape}((v.\text{shape}[0] * v.\text{shape}[1], v.\text{shape}[2]))$$

If we have "image" and want to return "v" as vector of shape ((length · height · 3), 1)

$$\hookrightarrow v = \text{image}.\text{reshape}((\text{image}.\text{shape}[0] * \text{image}.\text{shape}[1] * \text{image}.\text{shape}[2]))$$

From optimal Numpy Lab

Important

- np.exp(x) works for any np.array and applies e^x to each element
- $\delta(z) = \frac{1}{1+e^{-z}}$ $\delta'(z) = \delta(z) \cdot (1 - \delta(z))$
- image2vector is commonly used in deep learning
- np.shape is widely used, you'll see the keeping vector/matrix dimensions straight will go toward eliminating a lot of bugs
- numpy has efficient built-in functions
- broadcasting is extremely useful

Vectorization is very important in DL. It provides computational efficiency and clarity

$$\begin{aligned}
 & L_1(y, \hat{y}) = \sum_{i=0}^m |y^{(i)} - \hat{y}^{(i)}| \rightarrow \text{np.sum}(\cancel{\text{abs}}((y - \hat{y}).\text{abs})), L_2(y, \hat{y}) = \sum_{i=0}^m (y^{(i)} - \hat{y}^{(i)})^2 \\
 & \hookrightarrow \text{np.sum}(\text{up.dot}((y - \hat{y}).\text{abs}), (y - \hat{y}).\text{abs}))
 \end{aligned}$$

DL Spec [Week 2] course I

Coding Assignment - Overview

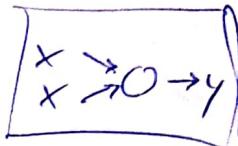
given:

training set $\rightarrow m_{\text{train}} \text{ images} \rightarrow \text{cat } (y=1)$

test set $\rightarrow m_{\text{test}} \text{ images} \rightarrow \text{non-cat } (y=0)$

Each image is of shape (num_px, num_px, 3)

each image is square $\frac{\text{num_px}}{\text{num_px}} \rightarrow (64, 64, 3) \rightarrow 12288$ $\stackrel{?}{\text{RGB}}$



\rightarrow train_set_x_orig, test_set_x_orig \rightarrow will be pre-processed \rightarrow train_set_x, test_set_x
 \rightarrow train_set_y = load_dataset()
 \rightarrow test_set_y

1) Find values for $m_{\text{train}} = \text{train_set_x_orig.shape}[0] \rightarrow 209$

~~figure out dimensions~~ $m_{\text{test}} = \text{test_set_x_orig.shape}[0] \rightarrow 50$

\rightarrow num_px = train_set_x_orig.shape[1] $\rightarrow 64$

2) Reshape training and test data sets into flattened arrays

with $X_{\text{flatten}} = X_{\text{original}}.reshape(X_{\text{original}}.shape[0], -1).T$

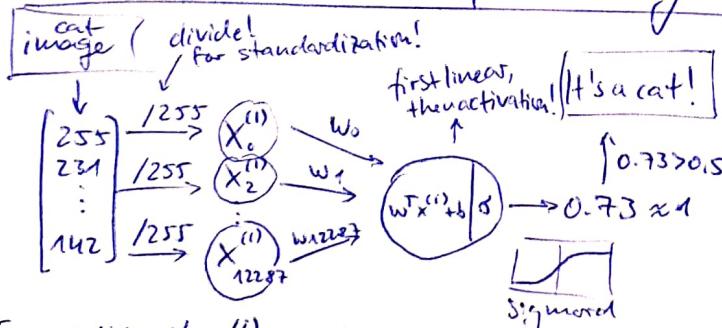
train_set_x扁平化 = train_set_x_orig.reshape(train_set_x_orig.shape[0], -1)
 " same for test_set_x_orig --"

3) Center and standardize dataset $\rightarrow (X - \mu)^{\text{mean}}$

transpose!

for pictures: divide every row by $\frac{\text{mean}}{\sigma - \text{std.deviation}}$
 by 255 (max. pixel value)

General Architecture of learning algorithm



For one example $x^{(i)}$:

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

$$a^{(i)} = \sigma(z^{(i)}) = \text{sigmoid}(z^{(i)})$$

$$\mathcal{L}(a^{(i)}, y^{(i)}) = -y^{(i)} \log(a^{(i)}) - (1-y^{(i)}) \log(1-a^{(i)})$$

The cost is computed by summing over all training examples

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

Key steps in this exercise

- Initialize all params
- Learn them by minimum cost
- Use learned params for prediction
- Analyze results and conclude

Main steps for building a Neural Network

1. Define model structure (e.g. # of input features)
2. Initialize the model's params
3. Loop:
 - 1 calculate current loss (FA)
 - " current gradient (BP)
 - Update parameters (Gradient Descent)

↳ build 1 → 3 separately and integrate them into one function model()

ex: create sigmoid as helper function

Initialize params

$$w = np.zeros((dim, 1))$$

for a picture (num_px * num_px * 3, 1)
= 64 * 64 * 3 = 12288

Forward propagation

- get X
- compute $A = \sigma(w^T X + b) = (a^{(1)} \dots a^{(m)})$ (logistic regression)
- cost $J = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)})$
- $\frac{\partial J}{\partial w} = \frac{1}{m} X(A-Y)^T$ $A = \text{sigmoid}(w^T X + b)$
 $\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$

$$X = \begin{bmatrix} 1 & 2 & -1 \\ 3 & 4 & -3.2 \end{bmatrix} \quad w = [1, 2]$$

$$Y = [1, 0, 1] \quad b = 2$$

Optimization: