

2 - Overview of the Problem set

Problem Statement: You are given a dataset ("data.h5") containing:

- a training set of `m_train` images labeled as cat (`y=1`) or non-cat (`y=0`)
- a test set of `m_test` images labeled as cat or non-cat
- each image is of shape `(num_px, num_px, 3)` where 3 is for the 3 channels (RGB). Thus, each image is square (`height = num_px`) and (`width = num_px`).

You will build a simple image-recognition algorithm that can correctly classify pictures as cat or non-cat.

Let's get more familiar with the dataset. Load the data by running the following code.

$(\text{num_px}, \text{num_px}, 3) \Rightarrow 3 \text{ channels (RGB)}$
square images

Exercise 2

Reshape the training and test data sets so that images of size `(num_px, num_px, 3)` are flattened into single vectors of shape `(num_px * num_px * 3, 1)`.

A trick when you want to flatten a matrix `X` of shape `(a,b,c,d)` to a matrix `X_flatten` of shape `(b*c*d, a)` is to use:

`X_flatten = X.reshape(X.shape[0], -1).T` # `X.T` is the transpose of `X`

```
In [6]: # Reshape the training and test examples
#(= 2 lines of code)
# train_set_x_flatten = ...
# test_set_x_flatten = ...
# YOUR CODE STARTS HERE
train_set_x_flatten = train_set_x_orig.reshape(train_set_x_orig.shape[0], -1).T
test_set_x_flatten = test_set_x_orig.reshape(test_set_x_orig.shape[0], -1).T

# YOUR CODE ENDS HERE

# Check that the first 10 pixels of the second image are in the correct place
assert np.alltrue(train_set_x_flatten[0:10, 1] == [196, 192, 190, 193, 186, 182, 188, 179, 174, 213]), "Wrong solution"
assert np.alltrue(test_set_x_flatten[0:10, 1] == [115, 110, 111, 137, 129, 129, 155, 146, 145, 159]), "Wrong solution"

print ("train_set_x_flatten shape: " + str(train_set_x_flatten.shape))
print ("train_set_y shape: " + str(train_set_y.shape))
print ("test_set_x_flatten shape: " + str(test_set_x_flatten.shape))
print ("test_set_y shape: " + str(test_set_y.shape))

train_set_x_flatten shape: (12288, 209)
train_set_y shape: (1, 209)
test_set_x_flatten shape: (12288, 50)
test_set_y shape: (1, 50)
```

$\text{train_set_x_flatten} \Rightarrow (m, \text{num_px}, \text{num_px}, 3)$
↳ RGB

$\text{reshape}(m, -1).T$

↳ automatically calculates product of the

Remaining dimensions.
($p_x \times p_x \times 3$)

$\frac{\cdot T}{\cdot 3}$ (feature examples).

now cd.

$$\begin{array}{l} \text{ex}_1 \downarrow \\ \text{ex}_2 \downarrow \\ \text{ex}_3 \downarrow \end{array} \begin{array}{c} \text{feature}_1 \rightarrow \\ \text{feature}_2 \rightarrow \\ \text{feature}_3 \rightarrow \\ \text{feature}_4 \rightarrow \end{array} \begin{bmatrix} 1 & 1 & 1 & 1 \\ x_1 & x_2 & x_3 & x_4 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

Preprocessing steps: \rightarrow figure out dimensions + shapes ($m, p_x, p_x, \text{batch}$)

\rightarrow reshape dataset 3, T.

each sample is ($p_x \times p_x \times 3, 1$)

\downarrow

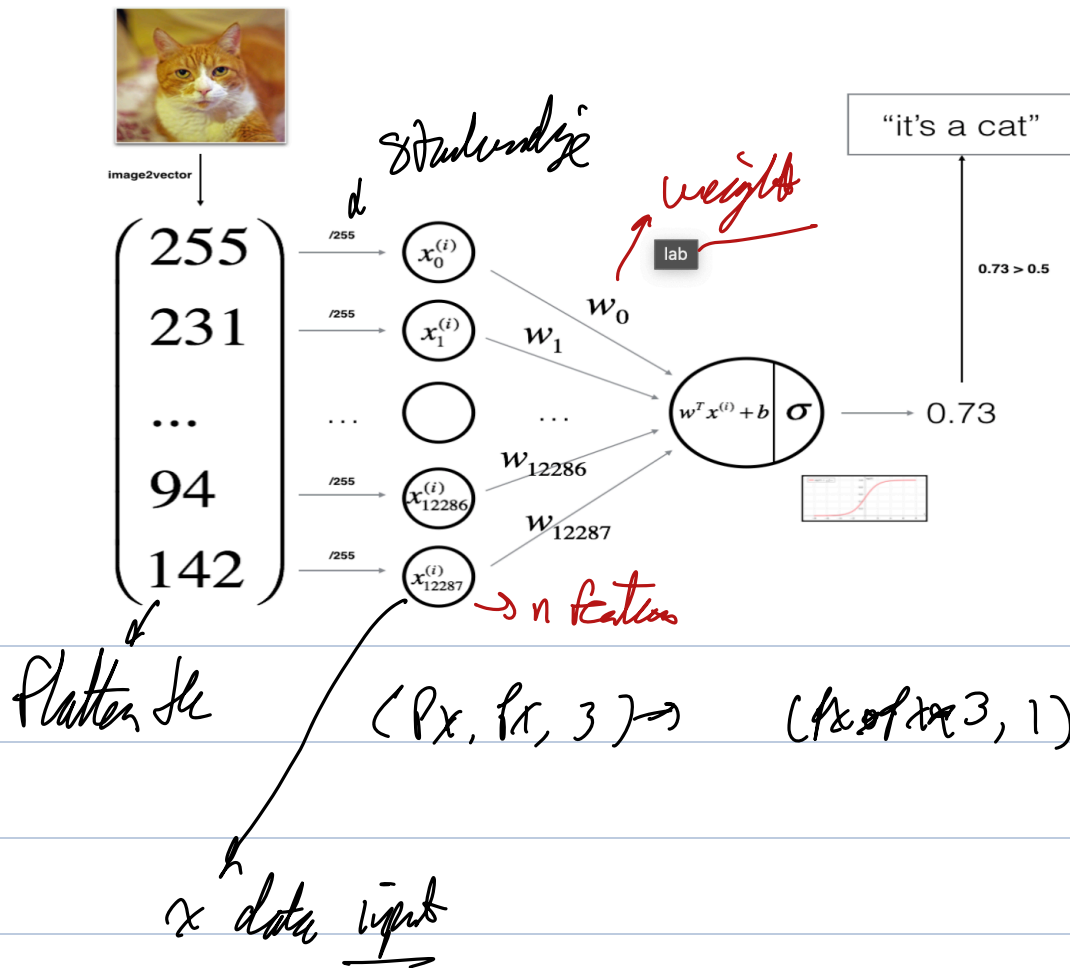
cd vector

\rightarrow Standardize \rightarrow divide by mean, or here divide by max value 255.

3 - General Architecture of the learning algorithm

It's time to design a simple algorithm to distinguish cat images from non-cat images.

You will build a Logistic Regression, using a Neural Network mindset. The following Figure explains why **Logistic Regression is actually a very simple Neural Network!**



$$y = \sigma(w^T x^{(i)} + b) \in \mathbb{R}$$

Handwritten notes:

- σ is the sigmoid function.
- w^T is the weight vector, size $1 \times n$.
- $x^{(i)}$ is the input vector, size $n \times 1$.
- b is the bias, size 1×1 .

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{n \times 1} \Rightarrow w^T = \begin{bmatrix} \vdots & \vdots & \vdots \end{bmatrix}_{1 \times n}$$

$y = 1$ for cat

$y = 0$ for not cat.

\rightarrow for example $x^{(i)}$

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

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

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

\downarrow
loss.

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

\rightarrow init the parameters of the model

\rightarrow Learn parameters for model by minimize cost.

\rightarrow use parameters to predict

\rightarrow analyze results

$$w = \underset{\text{features}}{\underbrace{\begin{bmatrix} 1 \\ \vdots \end{bmatrix}}_n} \quad x = \underset{\text{features}}{\underbrace{\begin{bmatrix} 1 & x_1 & x_2 \\ \vdots & \vdots & \vdots \end{bmatrix}}_m}$$

$$A = \sigma(w^T x + b)$$

np.dot(w, x)

$$\begin{bmatrix} w_0 & w_1 & w_2 & w_3 & w_4 \end{bmatrix}_{1 \times n}$$

$$\begin{bmatrix} \quad \quad \quad \end{bmatrix}_{1 \times m} + b.$$

$$\rightarrow 1 \times m.$$

$$A = 1 \times m.$$

no. examples.

$$\text{Cost} = \left(\frac{1}{m} \right) \sum_{i=1}^m y \log(A) + (1-y) \log(1-A)$$

$$1 \times m \quad 1 \times m$$

= element wise multiplication

Sum over them and then divide.

lab

```
# FORWARD PROPAGATION (FROM X TO COST)
#(≈ 2 lines of code)
# compute activation
# A = ...
# compute cost by using np.dot to perform multiplication.
# And don't use loops for the sum.
# cost = ...
# YOUR CODE STARTS HERE
A = sigmoid(np.dot(w.T,X) + b)
cost = (-1/m)*np.sum(Y*np.log(A) + (1-Y)*np.log(1-A))
```

$$w = \begin{bmatrix} \vdots \end{bmatrix} \quad \langle n, 1 \rangle \Rightarrow w^T = \begin{bmatrix} \quad \end{bmatrix} \langle 1, n \rangle$$

$$X = \begin{bmatrix} | & | & \dots & | \\ x_1 & x_2 & \dots & x_m \\ | & | & \dots & | \end{bmatrix} \quad \langle n, m \rangle$$

$$w^T X = \langle 1, n \rangle \langle n, m \rangle \Rightarrow \langle 1, m \rangle$$

b scalar \Rightarrow broadcasted to $\langle 1, m \rangle$

$$\Rightarrow A : \langle 1, m \rangle \quad \begin{bmatrix} \quad \end{bmatrix}$$

\hookrightarrow prediction for each feature in the m samples,

\hookrightarrow apply appropriate weight to appropriate input's feature.

+ which is the bias vectorized

$$\text{Cost} = \left(-\frac{1}{m} \right).$$

$$Y = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^{1, n} \quad \log(A) \Rightarrow \in \mathbb{R}^{1, n} \quad [\quad]$$

$Y \neq \log(A) \Rightarrow \in \mathbb{R}^{1, n}$ least vector multiplication.

Sum over the $\in \mathbb{R}^{1, n} \Rightarrow$ scalar cost
which is the -ve log likelihood cost.

b/c we want to minimize by making our
likelihood max \Rightarrow that we want
maximization.

Derivatives for the gradient descent

$$\frac{\partial J}{\partial W} = \frac{1}{n} X (A - Y)^T$$

\downarrow
here for a new rule to change.

$$\frac{\partial J}{\partial b} = \frac{1}{n} \sum (a^{(i)} - y^{(i)})$$

$$A \in \mathbb{R}^{1, n}$$

$$Y \in \mathbb{R}^{1, n} \quad \rightarrow^T \quad \in \mathbb{R}^{n, 1}$$

$$X \in \mathbb{R}^{n, n}$$

$$X \cdot (A - Y)^T = \in \mathbb{R}^{n, 1} \quad \uparrow \quad \text{Use every one of the } n \text{ weights}$$

$\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$

now much to

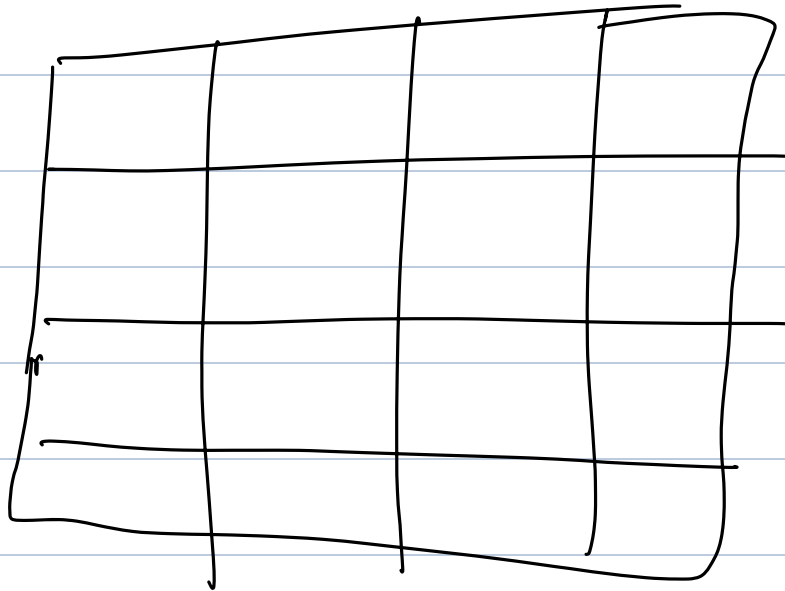
$$\frac{\partial T}{\partial b} = \frac{1}{n} \sum (a^{(i)} - y^{(i)}) + \frac{1}{n} \text{np.sum}(A, 1)$$

↪ how the bias

$\langle 1, n \rangle, \langle 1, n \rangle$.

Now we can compute cost and the gradient.

↪ want to update parameters using gradient descent.



Optimization For W we can try to minimize J.
Cost.

For $\theta = \theta - \alpha \text{d}\theta$
↪ learning rate.

$$w \Rightarrow \begin{bmatrix} w_0 \\ \vdots \\ w_n \end{bmatrix} \langle p_x \otimes p_y \otimes 3, 1 \rangle$$

$$n = p_x \otimes p_y \otimes 3$$

$$m = \# \text{ of } \underline{\text{expts}}$$

$$X = \begin{bmatrix} \begin{matrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{matrix} \end{bmatrix} \begin{matrix} 1 \\ n \\ 1 \end{matrix} \quad \begin{matrix} \xleftarrow{m} & \xrightarrow{m} \end{matrix} \quad \langle n, m \rangle$$

Return w , at b ,
 grads \rightarrow