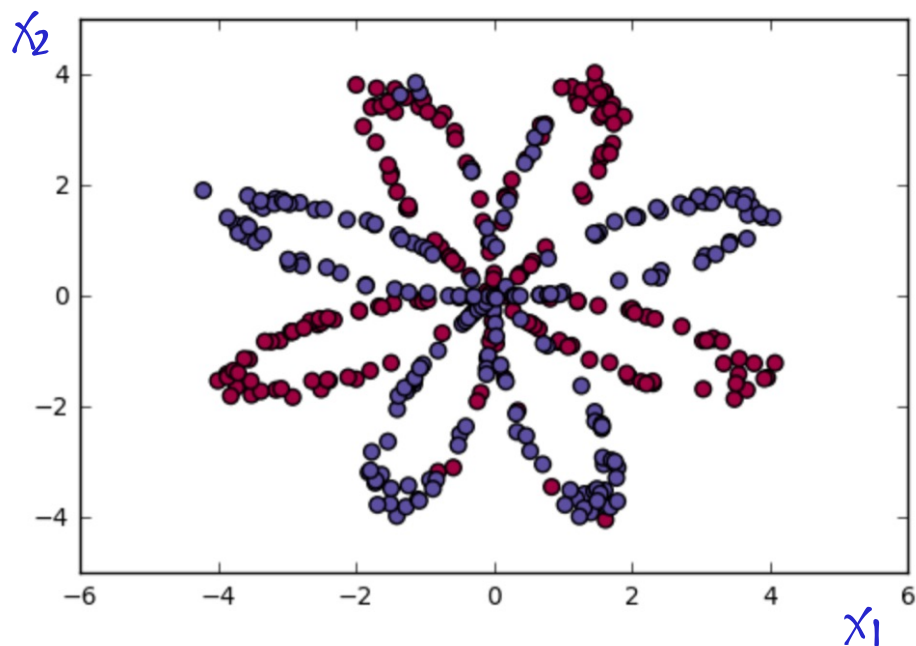# § planar data classification with a hidden layer

## 1. Dataset
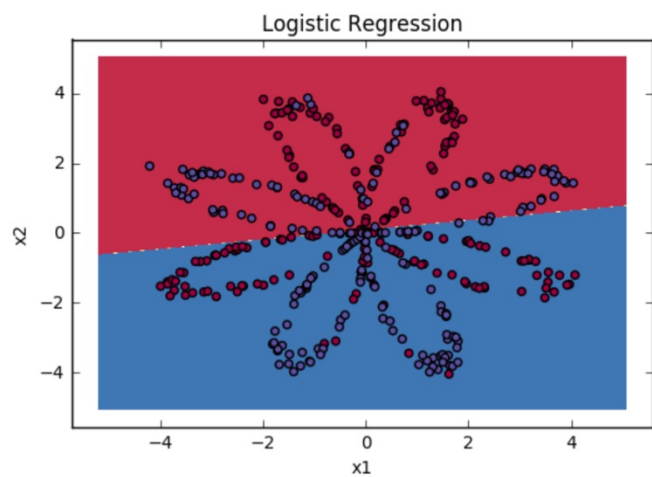
$X, Y =$ load_planar_dataset( )

$X$ — a numpy-array contains your features $(X_1, X_2)$
2×400

$Y$ — a numpy-array contains your Zables (red: 0, blue: 1)
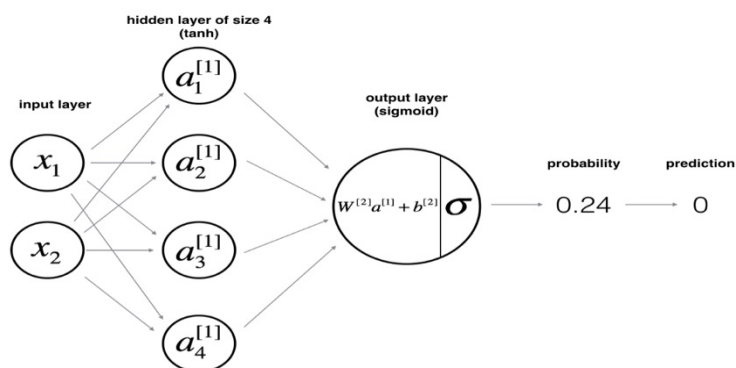1×400



## 2. Logistic Regression (Sample)



**Expected Output:**

| Accuracy | 47% |
|---|---|

The logistic regression doesn't perform well since the dataset is not linearly separable.

## 3. Neural Network model



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

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1](i)}$$
$$a^{[1](i)} = \tanh(z^{[1](i)})$$
$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2](i)}$$
$$y^{(i)} = a^{[2](i)} = \sigma(z^{[2](i)})$$
$$y^{(i)}_{prediction} = \begin{cases} 1 & \text{if } a^{[2](i)} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Given the predictions on all the examples, you can also compute the cost $J$ as follows:

$$J = -\frac{1}{m}\sum_{i=0}^{m}\left(y^{(i)}\log(a^{[2](i)}) + (1-y^{(i)})\log(1-a^{[2](i)})\right)$$

Reminder: The general methodology to build a Neural Network is to:
1. Define the neural network structure ( # of input units,  # of hidden units, etc).
2. Initialize the model's parameters
3. Loop:
    - Implement forward propagation
    - Compute loss
    - Implement backward propagation to get the gradients
    - Update parameters (gradient descent)

(1) define the neural network structure

$\begin{cases} n\_x : \text{the size of the input layer} \\ n\_h : \text{the size of the hidden layer} \\ n\_y: \text{the size of the output layer} \end{cases}$

(2) initialize the model's parameters
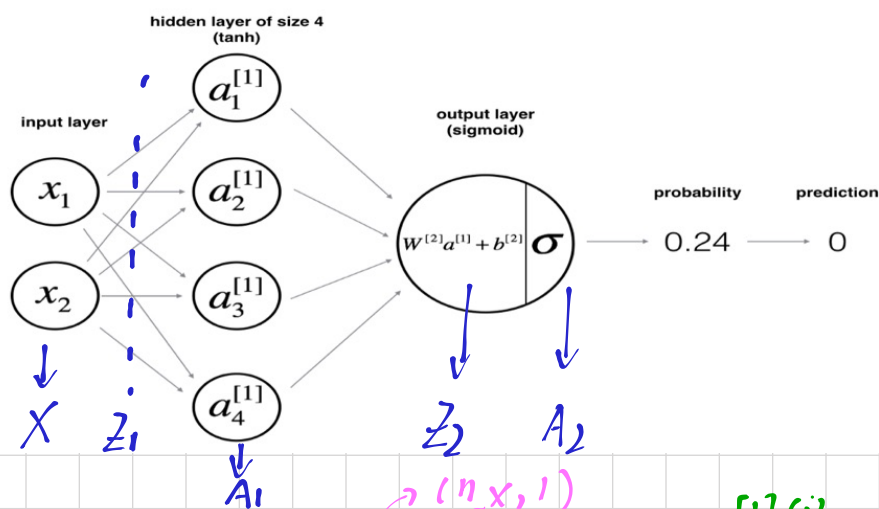
$W1 = np.random.randn (n\_h, n\_x) * 0.01 \rightarrow$ weight matrix

$b1 = np.zeros ((n\_h, 1)) \rightarrow$ bias vector

$W2 = np.random.randn (n\_y, n\_h) * 0.01 \rightarrow$ weight matrix

$b2 = np.zeros ( (n\_y, 1) ) \rightarrow$ bias vector

(3) loop



hidden layer of size 4
(tanh)

input layer

$a_1^{[1]}$
$a_2^{[1]}$
$a_3^{[1]}$
$a_4^{[1]}$

$x_1$
$x_2$

output layer
(sigmoid)

$W^{[2]}a^{[1]} + b^{[2]}$ $\sigma$

probability    prediction

0.24           0

X   $Z_1$

$A_1$

$Z_2$   $A_2$

$Z1 = np.dot (W1, X) + b1$ $\quad \overset{(n\_x, 1)}{\longrightarrow}$ $Z^{[1](i)} = W^{[1]} X^{(i)} + b^{[1](i)}$
$\quad (n\_h, n\_x) \quad (n\_h, 1)$

$A1 = np.tanh (Z1) \quad \longrightarrow a^{[1](i)} = \tanh (Z^{[1](i)})$
$\quad\quad (m\_h, 1)$

$Z2 = np.dot (W2, A1) + b2 \longrightarrow Z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2](i)}$
$\quad (n\_y, n\_h) \;(n\_h, 1) \quad \rightarrow (n\_y, 1)$

$A2 = sigmoid (Z2) \quad \longrightarrow a^{[2](i)} = \sigma (Z^{[2](i)})$
$\quad\quad (n\_y, 1)$

$$J = -\frac{1}{m} \sum_{i=0}^{m} \left( y^{(i)} \log\left(a^{[2](i)}\right) + (1 - y^{(i)}) \log\left(1 - a^{[2](i)}\right) \right)$$

$\rightarrow (n\_y, 1)$

$\{$ logprobs = np.multiply ( np.log (A2) , $\overset{(n\_y,1)}{Y}$ ) + np.multiply (np.log (1-A2) , $\overset{(n\_y,1)}{(1-Y)}$ )

$[$ cost = -1/m $\times$ np.sum (logprobs)

### P.S. Broadcasting

When operating on two arrays, Numpy compares their shapes element-wise.

Two dimensions are compatible when:

$\{$ ① they are equale , or     others → Value Error

$[$ ② one of them is /

e.g. ① A: 8 X①X⑥X① ②     ② A: ①⑤X③X⑤①     ③ A: 8 X④X3   → not 1/0

B: ⑦X①X⑤             B: ⑤X①X⑤             B: ②X1

result: 8 X7X6 X1     result: 15 X3 X5     result: Error

### Backward propagation

## Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y \qquad\qquad dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = dz^{[2]}a^{[1]T} \qquad\quad dW^{[2]} = \frac{1}{m}dZ^{[2]}A^{[1]T}$$

$$db^{[2]} = dz^{[2]} \qquad\qquad\quad db^{[2]} = \frac{1}{m}np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

$$dz^{[1]} = W^{[2]T}dz^{[2]} * g^{[1]\prime}(z^{[1]}) \quad dZ^{[1]} = W^{[2]T}dZ^{[2]} * g^{[1]\prime}(Z^{[1]})$$

$$dW^{[1]} = dz^{[1]}x^{T} \qquad\qquad dW^{[1]} = \frac{1}{m}dZ^{[1]}X^{T}$$

$$db^{[1]} = dz^{[1]} \qquad\qquad\quad db^{[1]} = \frac{1}{m}np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

Andrew Ng

$\rightarrow (n\_y, 1)$

dZ2 = A2 - Y

$\rightarrow (n\_y, n\_h)$    $\rightarrow (n\_y,1)$    $\rightarrow (1, n\_h)$

dW2 = (1/m) $\times$ np.dot (dZ2 , A1. T)

$\rightarrow (n\_y, 1)$

db2 = (1/m) $\times$ np.sum (dZ2 , axis=1 , keepdims = True)

$\rightarrow (n\_h, 1)$    $\rightarrow (n\_h, 1)$    $\rightarrow (n\_h, 1)$

dZ1 = np.multiply( np.dot (W2.T, dZ2) , (1- np.power (A1,2)))

$(n\_h, n\_y)$    $(n\_y, 1)$

dW1 = (1/m) $\times$ np.dot (dZ1, X.T)

$(n\_h, n\_X)$                    $(n\_h, 1)$    $(1, n\_X)$

db1 = (1/m) $\times$ np.sum (dZ1, axis=1, keepdims=True)

$(n\_h, 1)$

$$dZ^{[1]} = W^{[2]^T} dZ^{[2]} \times g^{[1]'} (Z^{[1]})$$

dZ1 = np.multiply (np.dot (W2.T, dZ2) , (1-np.power (A1,2)))

(n_h, n_y)  (n_y, 1) same          (n_h, 1)

(n_h, 1)

Broadcasting

(n_h, 1)

✗ attention the difference
etween np.dot() and
np.multiply

**Gradient rule**

$$\theta = \theta - \alpha \frac{\partial J}{\partial \theta}$$

(4) Integrate parts in nn_model

① Initialize parameters

parameters = initialize_parameters (n_x, n_h, n_y)

② Loop

for i in range (0, num_iterations):

   i. Forward propagation

   A2, cache = forward_propagation (X, parameters)

   ii. Cost Function

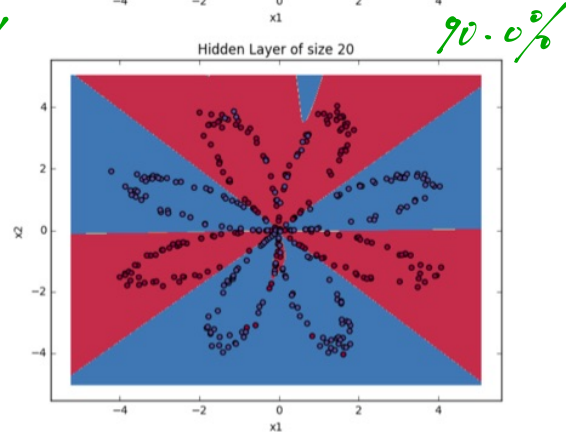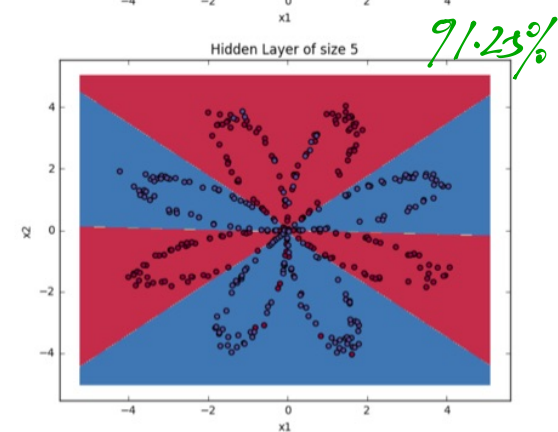   cost = compute_cost (A2, Y, parameters)

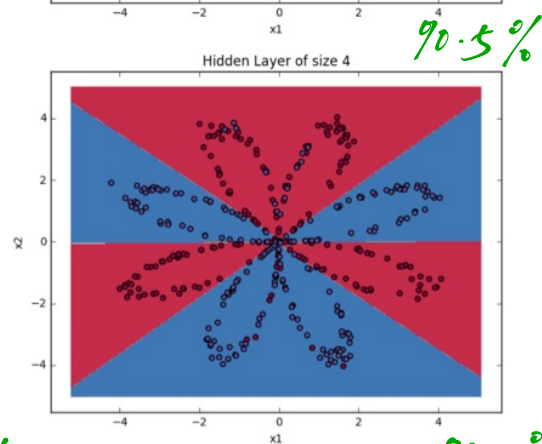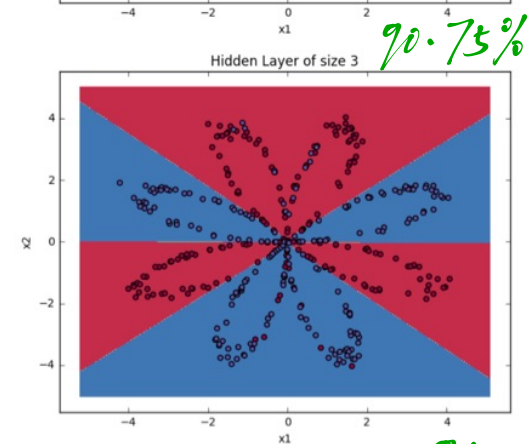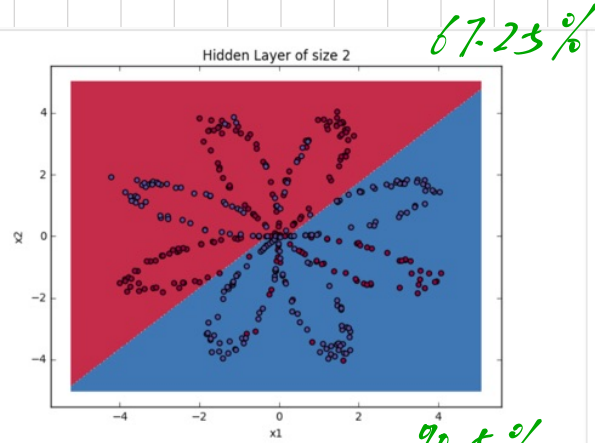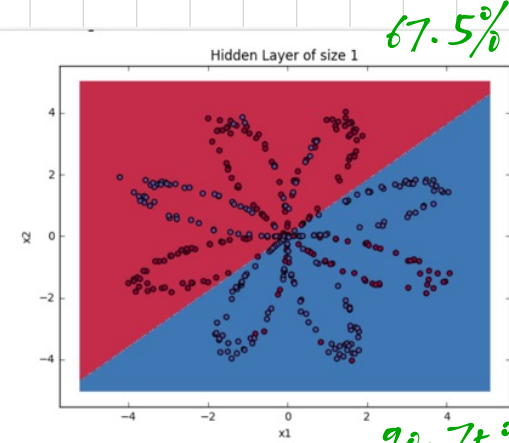   iii. Back propagation

   grads = backward_propagation (parameters, cache, X, Y)

   iv. Parameters update

   parameters = update_parameters (parameters, grads)

(5) Turning hidden layer size

change the hidden layer size from 1 → 2 → 3 → 4 → 5 → 20 → 50

Hidden Layer of size 1 — 67.5%

Hidden Layer of size 2 — 67.25%

Hidden Layer of size 3 — 90.75%

Hidden Layer of size 4 — 90.5%

Hidden Layer of size 5 — 91.25%

Hidden Layer of size 20 — 90.0%

Hidden Layer of size 50 — 90.25%

Interpretation:
The larger models (with more hidden units) are able to fit the training set better, until eventually the largest models overfit the data.
The best hidden layer size seems to be around n_h = 5. Indeed, a value around here seems to fits the data well without also incurring noticable overfitting.
You will also learn later about regularization, which lets you use very large models (such as n_h = 50) without much overfitting.