

Neural Network Basics

1. how to convert the image to vector:

Binary Classification

In a binary classification problem, the result is a discrete value output.

For example

- account hacked (1) or not hacked (0)
- a tumor malign (1) or benign(0)

Example: Cat vs Non-Cat

The goal is to train a classifier for which the input is an image represented by a feature vector, x , and predicts whether the corresponding label y is 1 or 0. In this case, whether this is a cat image(1) or a non-cat image (0).



An image is stored in the computer in three separate matrices corresponding to the Red, Green, and Blue color channels of the image. The three matrices have the same size as the image, for example, the resolution of the cat image is 64 pixels X 64 pixels, the three matrices (RGB) are 64 X 64 each.

The value in a cell represents the pixel intensity which will be used to create a feature vector of n-dimension. In pattern recognition and machine learning, a feature vector represents an image. Then the classifier's job is to determine whether it contains a picture of a cat or not.

To create a feature vector, x , the pixel intensity values will be “unrolled” or “reshaped” for each color. The dimension of the input feature vector x is $n = 64 \times 64 \times 3 = 12288$.

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix} \left\{ \begin{array}{l} \text{red} \\ \text{green} \\ \text{blue} \end{array} \right.$$

Notation

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

m training examples : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$M = M_{\text{train}}$$

$$M_{\text{test}} = \# \text{test examples.}$$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}^{\uparrow} \quad X \in \mathbb{R}^{n_x \times m} \quad X \cdot \text{shape} = (n_x, m)$$

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

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

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

But to take the stuff or

当有 m 个training instance时，将 m 个instance横向堆叠所以feature vector X 的dimension就是 $n_x \times m$, Y 也横向堆叠, $Y \cdot \text{shape} = (1, m)$

2. logistic regression:

Logistic Regression

Logistic regression is a learning algorithm used in a supervised learning problem when the output y are all either zero or one. The goal of logistic regression is to minimize the error between its predictions and training data.

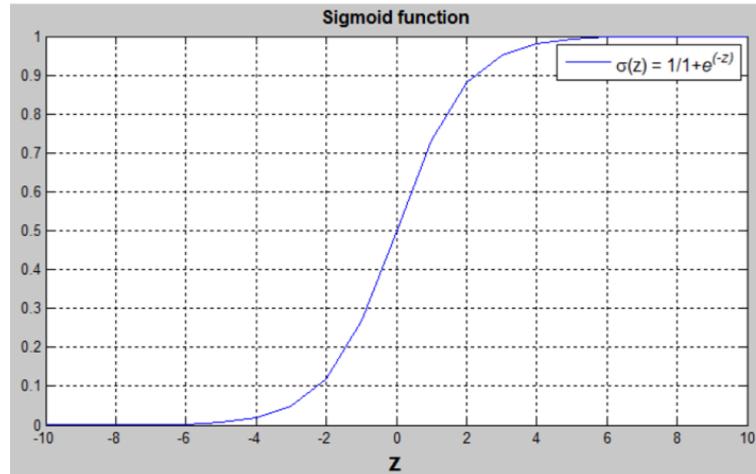
Example: Cat vs No - cat

Given an image represented by a feature vector x , the algorithm will evaluate the probability of a cat being in that image.

$$\text{Given } x, \hat{y} = P(y = 1|x), \text{ where } 0 \leq \hat{y} \leq 1$$

The parameters used in Logistic regression are:

- The input features vector: $x \in \mathbb{R}^{n_x}$, where n_x is the number of features
- The training label: $y \in \{0,1\}$
- The weights: $w \in \mathbb{R}^{n_x}$, where n_x is the number of features
- The threshold: $b \in \mathbb{R}$
- The output: $\hat{y} = \sigma(w^T x + b)$
- Sigmoid function: $s = \sigma(w^T x + b) = \sigma(z) = \frac{1}{1+e^{-z}}$



$(w^T x + b)$ is a linear function ($ax + b$), but since we are looking for a probability constraint between $[0,1]$, the sigmoid function is used. The function is bounded between $[0,1]$ as shown in the graph above.

Some observations from the graph:

- If z is a large positive number, then $\sigma(z) = 1$
- If z is small or large negative number, then $\sigma(z) = 0$
- If $z = 0$, then $\sigma(z) = 0.5$

3. The 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.

Logistic Regression: Cost Function

To train the parameters w and b , we need to define a cost function.

Recap:

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

$x^{(i)}$ the i-th training example

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, we want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (error) function:

The loss function measures the discrepancy between the prediction ($\hat{y}^{(i)}$) and the desired output ($y^{(i)}$). In other words, the loss function computes the error for a single training example.

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2$$

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

- If $y^{(i)} = 1$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$ where $\log(\hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 1
- If $y^{(i)} = 0$: $L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$ where $\log(1 - \hat{y}^{(i)})$ and $\hat{y}^{(i)}$ should be close to 0

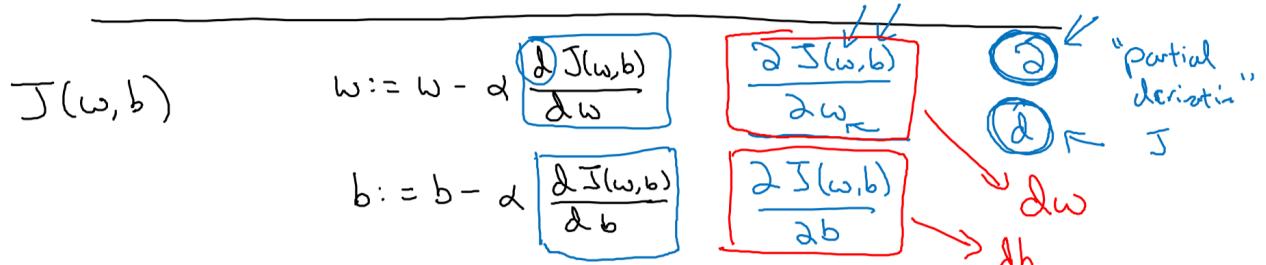
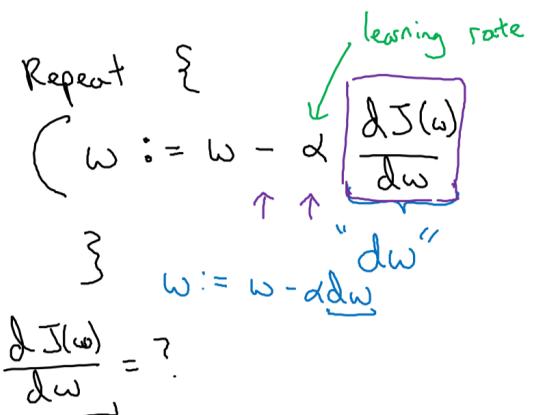
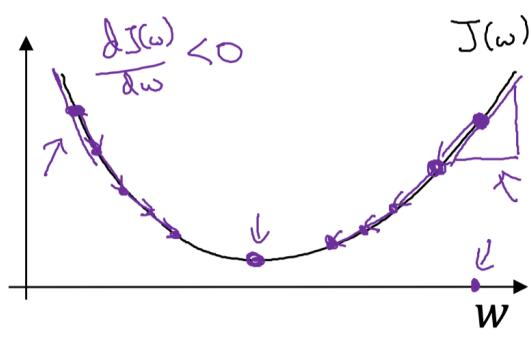
Cost function

The cost function is the average of the loss function of the entire training set. We are going to find the parameters w and b that minimize the overall cost function.

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

4. Gradient Descent: w 和**b**更新方式, 这里alpha是learning rate也就是步长

Gradient Descent



5. computation graph: 从左到右是forward propagation, 从右到左是backwards propagation,

Computing derivatives

Handwritten notes for computing derivatives:

$$J = 3v$$

$$v = a + u$$

$$u = bc$$

$$a = 5 \quad \frac{\partial J}{\partial a} = 3$$

$$b = 3 \quad \frac{\partial J}{\partial b} = 6$$

$$c = 2 \quad \frac{\partial J}{\partial c} = 9$$

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

$$\frac{\partial J}{\partial a} = 3 = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial a} = 3 \cdot 1$$

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial b} = 3 \cdot 2 = 6$$

$$\frac{\partial J}{\partial c} = \frac{\partial J}{\partial u} \cdot \frac{\partial u}{\partial c} = 3 \cdot 3 = 9$$

$$u = 6 \rightarrow 6.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$b = 3 \rightarrow 3.001$$

$$u = b \cdot c = 6 \rightarrow 6.002$$

$$J = 33.006$$

$$v = 11.002$$

$$J = 33V$$

Andre

In the context of neural networks and machine learning, forward propagation and backward propagation are two fundamental processes that occur during the training of a neural network. They are essential for understanding how neural networks learn from data.

1. Forward Propagation:

- Definition:** Forward propagation is the process of moving the input data through the neural network to generate an output. This is essentially how the neural network makes predictions.
- Process:** During forward propagation, the input data is fed into the network. It then passes through the layers of the network, where each layer applies specific weights and biases, often followed by a non-linear activation function. The result is the predicted output at the end.
- Purpose:** The main purpose of forward propagation is to compute the output based on the current state of the network's weights and biases.

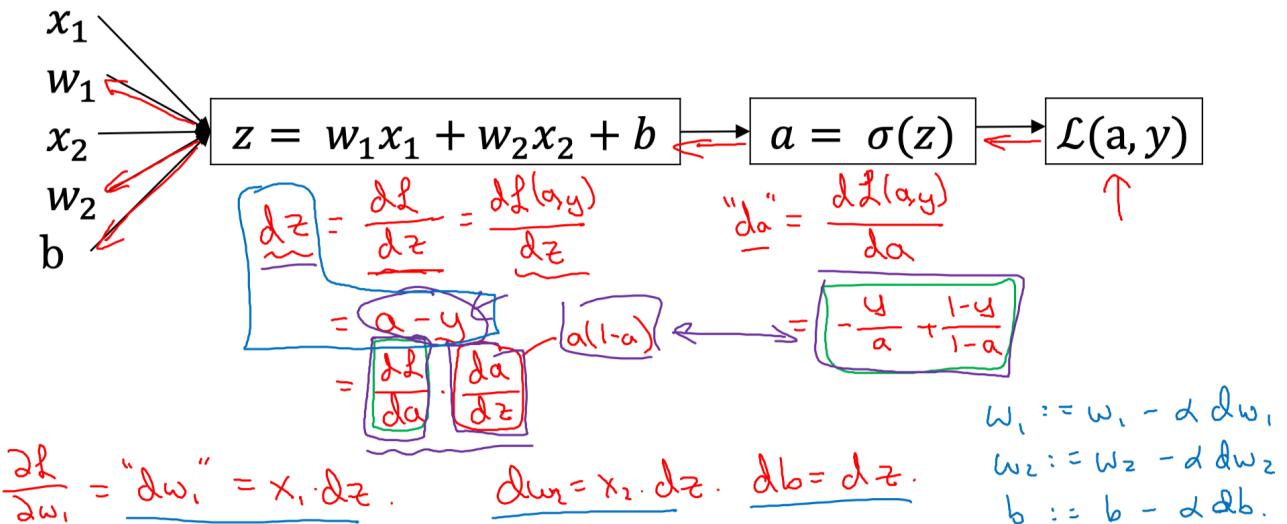
2. Backward Propagation:

- Definition:** Backward propagation, often called backpropagation, is the process of tuning the weights and biases of the neural network based on the error of the prediction. It is a key part of the training process.
- Process:** In backward propagation, the error (or loss) is calculated by comparing the predicted output of the network with the actual target values. This error is then propagated back through the network, starting from the output layer to the input layer. During this process, the weights and biases are adjusted in order to minimize the error.
- Purpose:** The main purpose of backward propagation is to optimize the weights and biases so that the network's predictions become more accurate over time.

In a typical training cycle of a neural network, forward propagation is used to calculate the output for a batch of data, and then backward propagation is used to update the network's weights and biases based on the error of that output. This cycle is repeated many times with different batches of data until the network becomes proficient at making accurate predictions.

6. logistic regression derivatives for one example:

Logistic regression derivatives



$$\begin{aligned}
 & \hat{y} \quad y \\
 & L(a, y) = -(y \log a + (1-y) \log(1-a)) \\
 & \boxed{\frac{\partial L}{\partial a}} = \frac{\partial L(a, y)}{\partial a} = \left[\frac{y}{a} + \frac{1-y}{1-a} \right] \\
 & \boxed{\frac{\partial L}{\partial z}} = \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \times \frac{\partial a}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = -y(1-a) + (1-y)a \\
 & \frac{\partial a}{\partial z} = \frac{e^{-z}}{1+e^{-z}} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} = a(1-a)
 \end{aligned}$$

反向传播链式求导得到L对于w1,w2,b的偏导,然后根据右下角更新这3个值

7. logistic regression derivatives for m example: 对6中取到的单个的derivative求和取平均

Logistic regression on m examples

$$\begin{aligned} J(\omega, b) &= \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)}) \\ \rightarrow a^{(i)} &= \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^\top x^{(i)} + b) \end{aligned}$$

$(x^{(i)}, y^{(i)})$
 $\underline{dw_1}, \underline{dw_2}, \underline{db}$

$$\underline{\frac{\partial}{\partial \omega_i} J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_i} l(a^{(i)}, y^{(i)})}_{\underline{dw_1^{(i)}} - (x^{(i)}, y^{(i)})}$$

$$\begin{aligned} J &= 0; \underline{dw_1} = 0; \underline{dw_2} = 0; \underline{db} = 0 \\ \rightarrow \text{For } i &= 1 \text{ to } m \\ z^{(i)} &= \omega^\top x^{(i)} + b \\ Q^{(i)} &= \sigma(z^{(i)}) \\ J &+= -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})] \\ \underline{dz^{(i)}} &= a^{(i)} - y^{(i)} \\ \begin{cases} \underline{dw_1} += x_1^{(i)} dz^{(i)} \\ \underline{dw_2} += x_2^{(i)} dz^{(i)} \\ \underline{db} += dz^{(i)} \end{cases} & n=2 \\ J / &= m \leftarrow \\ \underline{dw_1 /} &= m; \underline{dw_2 /} = m; \underline{db /} = m. \leftarrow \end{aligned}$$

$$\underline{dw_1} = \frac{\partial J}{\partial \omega_1}$$

$$\omega_1 := \omega_1 - \alpha \underline{dw_1}$$

$$\omega_2 := \omega_2 - \alpha \underline{dw_2}$$

$$b := b - \alpha \underline{db}.$$

Vectorization

Δn

这里需要2个loop，一个loop遍历 m 个instance，另一个loop遍历 n 个参数，但是嵌套的loop太耗时了，所以用vectorization

Python and Vectorization

1. python and vectorization:

```
import time
import numpy as np

a=np.random.rand(1000000)
b=np.random.rand(1000000)
tic=time.time()
c=np.dot(a,b)
toc=time.time()

print(c)
print("Vectorized version:"+str(1000*(toc-tic))+"ms")
```

```

c=0
tic=time.time()
for i in range(1000000):
    c+=a[i]*b[i]
toc=time.time()
print(c)
print("For loop:"+str(1000*(toc-tic))+"ms")

```

result:

```

/usr/bin/python3 /Users/nanzhenghan/PycharmProjects/Bilibili/main.py
250166.45798134204
Vectorized version:0.5202293395996094ms
250166.45798133797
For loop:187.50977516174316ms

```

Process finished with exit code 0

2. set dw as a vector:

Logistic regression derivatives

$J = 0, \boxed{dw_1 = 0, dw_2 = 0}, db = 0$ $dw = np.zeros((n_x, 1))$

\rightarrow for i = 1 to n:

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

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

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

$$\boxed{dz^{(i)} = a^{(i)}(1 - a^{(i)})}$$

\downarrow for j = 1 ... n_x

$$\boxed{dw_1 += x_1^{(i)} dz^{(i)}, dw_2 += x_2^{(i)} dz^{(i)}, db += dz^{(i)}} \quad | \quad n_x = 2$$

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

$$J = J/m, \boxed{dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m}$$

$$dw /= m.$$

3. vectorizing logistic regression:

Vectorizing Logistic Regression

$$\begin{aligned} \rightarrow z^{(1)} &= w^T x^{(1)} + b \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) \end{aligned}$$

$$z^{(2)} = w^T x^{(2)} + b$$

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

$$z^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

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

$\frac{(n_x, m)}{\mathbb{R}^{n_x \times m}}$

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

$\rightarrow \underline{z} = \begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = w^T \underline{X} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} w^T X^{(1)} + b \\ w^T X^{(2)} + b \\ \vdots \\ w^T X^{(m)} + b \end{bmatrix}_{1 \times m}$

$\rightarrow \underline{z} = np.\dot{e}(w.T, X) + b_{(1,1)}$ "Broadcasting"

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

Andrew Ng

可以用 $z = np.dot(w.T, x) + b$ 来直接算 z 来避免 for loop

4. Implementing logistic regression with complete vectorization (code on the right), but we can't get rid of the iteration number of gradient descent:

Implementing Logistic Regression

```
J = 0, dw1 = 0, dw2 = 0, db = 0
for i = 1 to m:
    z(i) = wTx(i) + b
    a(i) = σ(z(i))
    J += -[y(i) log a(i) + (1 - y(i)) log(1 - a(i))]
    dz(i) = a(i) - y(i)
    dw1 += x1(i) dz(i)
    dw2 += x2(i) dz(i)
    db += dz(i)
J = J/m, dw1 = dw1/m, dw2 = dw2/m
db = db/m
```

```
for iter in range(1000):
    z = wTX + b
    = np.dot(w.T, X) + b
    A = σ(z)
    dZ = A - Y
    dw = 1/m * X * dZT
    db = 1/m * np.sum(dZ)
    w := w - αdw
    b := b - αdb
```