

Week 2: Basics of NN programming.
→ Binary Classification

- When processing the training set we process the training set without necessarily using the for statement
- Logistic regression is an algorithm for binary class.

Converted [image] → 1 (Cat) vs 0 (non Cat)
to Feature vector (\mathbf{x}) → y .

• Notations:

Single training sample is denoted (\mathbf{x}, y) where $\mathbf{x} \in \mathbb{R}^{n_x}$

\mathbf{x} is a feature vector in dimension n_x
 y is 0 or 1

m → train examples.

$m = m_{train}$ $m_{test} = \#$ test examples

to put all of training examples in compact notation
we define matrix:

This matrix takes $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$
as stack them to columns.

$$\mathbf{X} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \mathbf{x}^{(3)} & \\ & & & \\ & & & \end{bmatrix}_{n_x \times m}$$

$$\mathbf{y} = \begin{bmatrix} y^{(1)} & \cdots & y^{(m)} \end{bmatrix}$$

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

→ logistic Regression

- This algorithm we use it when the output label y is $0 \text{ or } 1$.

given X , want $\hat{y} = P(y=1|X) \quad 0 \leq \hat{y} \leq 1$
 \rightarrow means the prob. that y is 1 given the input feature X

ex: parameters $w \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}$.

output $\hat{y} = \underbrace{\sigma(w^T x + b)}_{\text{sigmoid}} \quad z$

$$\sigma(z) = \frac{1}{1+e^{-z}} \rightarrow \begin{aligned} &\text{if } z \text{ is large then } \sigma(z) \\ &\text{will be near 1 then it} \\ &\frac{1}{1+0} = 1 \end{aligned}$$

$$\text{if } z \text{ is true } \approx \frac{1}{1+\text{num}} = 0$$

to train the parameters a, b we need to define cost fn.

→ logistic regression Cost Function

ex: MSE loss function: measures error for one training sample (x, y)
 - cost function: Avg. error for all training samples

loss or error fun: $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \rightarrow$ measures how good our output \hat{y} is while true output label y is

- in logistic regression we write loss fun:

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

& this lossfun. is called binary crossentropy.

$\therefore y=1 \ L(\hat{y}, y) = -\log \hat{y}$, so we want \hat{y} large

$\therefore y=0 \ L(\hat{y}, y) = -\log(1-\hat{y})$, so $\log(1-\hat{y})$ large so \hat{y} small

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

→ gradient Descent

• recap: $\hat{y} = \sigma(w^T x + b)$, $\sigma(z) = \frac{1}{1+e^{-z}}$ } logistic reg alg.

• Cost fun : $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$

• loss fun measures how well the alg. predicts y for each training sample compared to the true value

• So we always try to find w, b that minimize $J(w, b)$

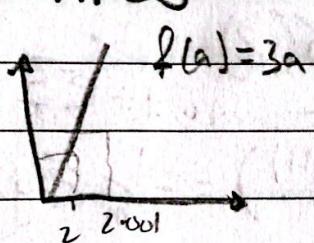
• gradient descent is optimization algorithm to tech. use to minimize the cost fun. by iteratively adjust the model parameter. So we improve model performance

$$\Theta = \Theta - \alpha \cdot \nabla J(\Theta)$$

$\Theta \rightarrow$ model parameters, α learning rate, $\nabla J(\Theta)$ gradient of cost func.

- . when calling a model we don't implement Cost or loss fun.
 - ① we create model ② it has its own built-in cost-fun.
ex: binary cross entropy for logistic reg. & its optimization algorithm (ex: gradient descent)
 - ③ when we call .fit(), the model calculate loss & cost fun.
update the model parameters
iterate this process until /
 - ④ prediction
- Convergence → when updates results no change in cost fun.

→ Derivatives:

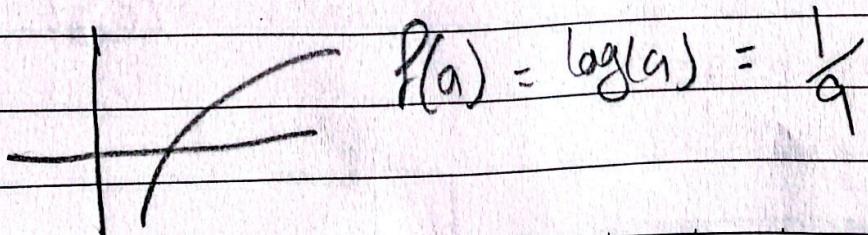
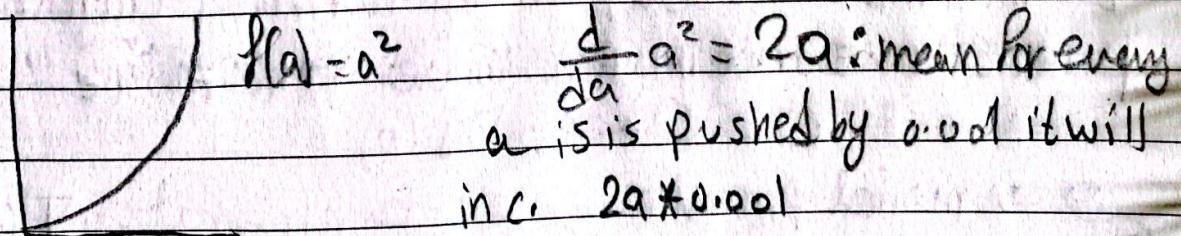


$$a = 2 \quad f(a) = 6 \\ a = 2.001 \quad f(a) = 6.003$$

So Slope (derivative) of $f(a)$
at $a = 2$ is 3

- one property of derivative whatever you take the slope of the fun. 3 whether $a = 2$ or 5 the slope = 3 where the inc. is 0.001

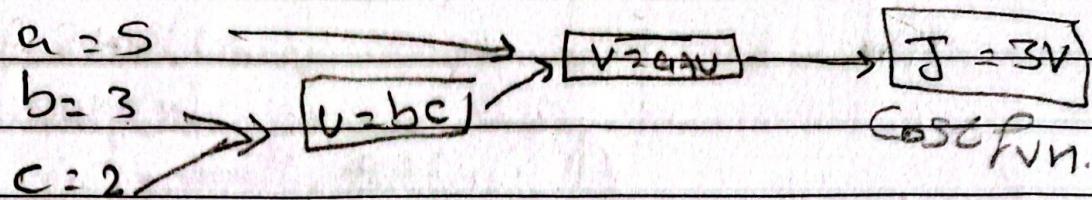
→ more Derivatives example.



→ Computation graph

$$J(a, b, c) = 3(a + bc) = 3(5 + 3 \times 2) = 33$$

$$U = bc, V = a + U, J = 3V$$



→ Derivatives with Computation graphs *

backpropagation

→ logistic regression Gradient descent *

. we will talk about how to compute derivatives to implement gradient descent

. logistic regression recap

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

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

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

. lets say we have two features : to compute z we put w_1, w_2 , we modify w & b to reduce this loss

x_1
 w_1

x_2

w_2

b

$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} \rightarrow L(a, y)$$

→ Vectorization.

• vectorization is the art of getting rid of loops in your code & this causes speed up your code.

• What is vect. ex: $z = w^t x + b$ $w = [i] \times [i]$

to compute $w^t x$ if we had a non vectorized implementation (non vectorized)

$$z = 0$$

for i in range(n-x)

$$z[i] = w[i] * x[i]$$

$$z += b$$

$$z = np.dot(w, x) + b$$

• GPU] both parallelization instruction or SIMD
CPU]

single instruction
multiple data

☞ More vectorization examples.

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

✓

$V = np.zeros((n, 1))$

for i in range(n):

$V[i] = math.exp(V[i])$

* Vectorizing logistic regression. ✓

☞ broad casting: is another tech. That make py code run faster

☞ Py & numpy vectors: $a = np.random.randn(5)$
 $\rightarrow a.shape = (5,) \rightarrow rank 1 array$
 $\rightarrow assert()$