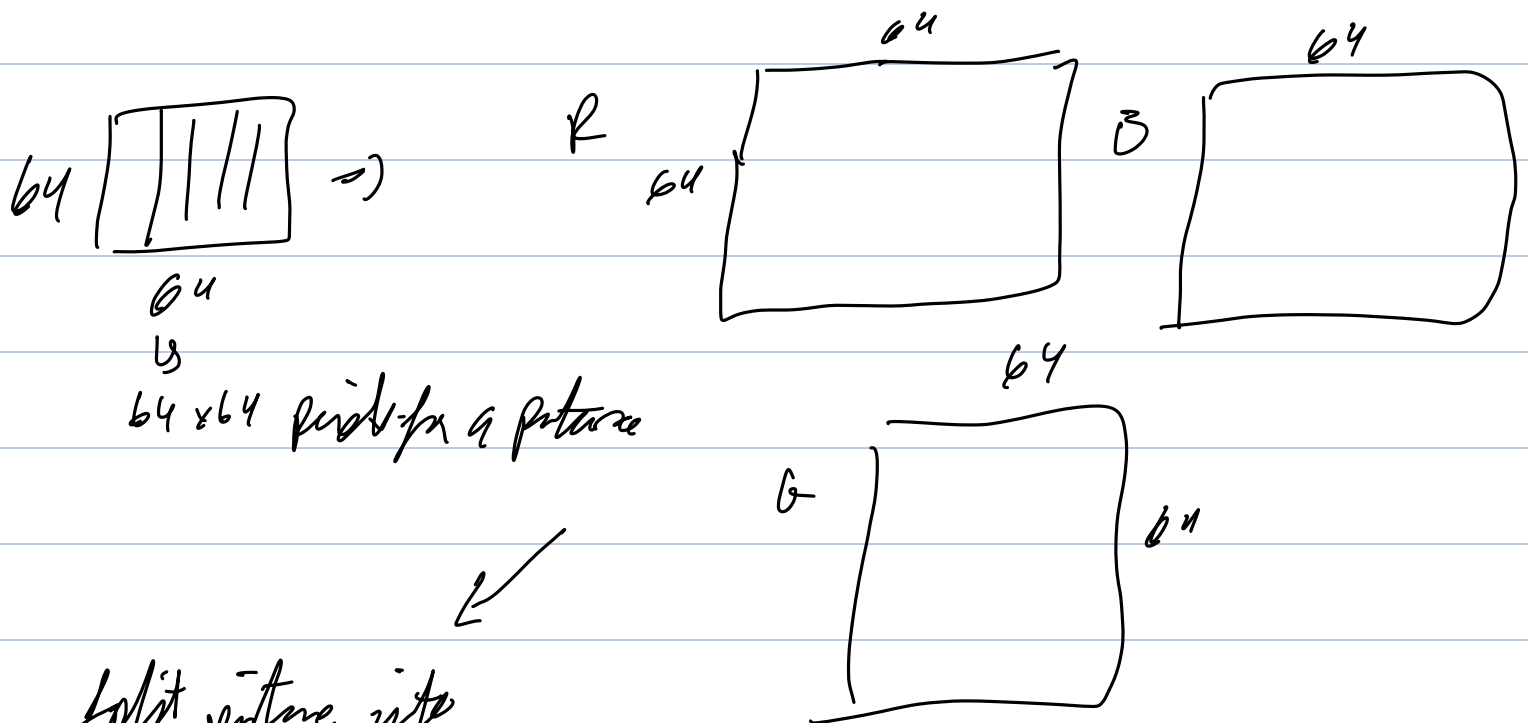


Logistic Regression as a Neural Network

Lec 1: Binary Classification

Logistic Regression for Binary Classification

Prin of Cat $\rightarrow y = 1$ (Cat) vs 0 (non cat)



Split picture into
the R, G, B matrices of
the picture

Single Feature Vector $\rightarrow X =$

$\begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 192 \\ \vdots \end{bmatrix}$	$\left. \begin{array}{l} \text{all red } 64 \times 64 \\ + \\ \text{all blue } 64 \times 64 \end{array} \right\}$
--	---

$\begin{bmatrix} \vdots \\ 142 \\ 172 \\ \vdots \end{bmatrix}$ } all green 64×64
 \downarrow
 x is a 12288 dim feature \Rightarrow $64 \times 64 \times 3 = 12288$
vector

$n = n_x = 12288 \rightarrow$ dimension of this input vector / feature

$x \rightarrow y (100)$
 \rightarrow binary classification

NOTATION

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

M training examples : $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, ..., $(x^{(M)}, y^{(M)})$

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

or Test set $M_{\text{test}} = \# \text{ test example}$

X
 \downarrow
 MATRIX

n_x

M

$\begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(M)} \\ | & | & & | \end{bmatrix}$

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

$$X_{\text{shape}} = (n_x, m) \quad \langle n_x, m \rangle$$

\hookrightarrow n_x by m dimensional matrix

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

\downarrow
matrix

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

$$Y_{\text{shape}} = (1, m) \quad \langle 1, m \rangle$$

Lec 2 Logistic Regression

Given x , want $\hat{y} = P(y=1 | x)$

if x is a picture, $\hat{y} =$ what is chance that x is a cat.

\hookrightarrow $1 \leq \hat{y} \leq 1$

$$x \in \mathbb{R}^{n_x}$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$.

weights bias

for logistic reg.

Output $\hat{y} = w^T x + b \rightarrow$ linear form of x

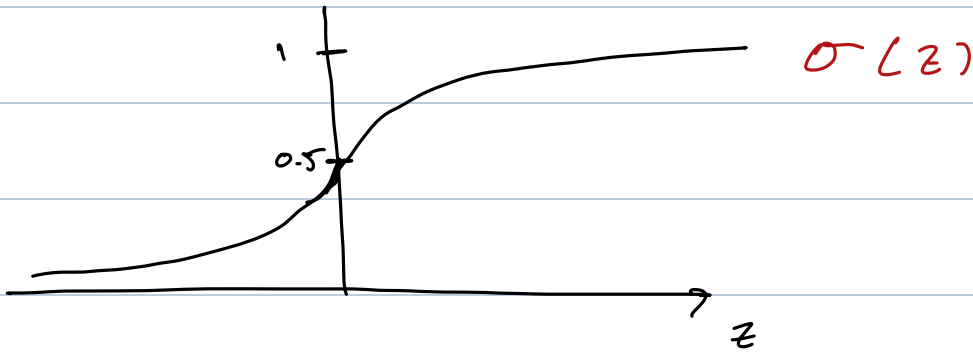
$w^T = 1 \times n_x$
 $x = n_x \times 1$
 $\hat{y} = 1 \times 1$

1×1

SOPAD (This is linear regression)
 \hookrightarrow not good b/c we want $\hat{y} \in [0, 1]$
which is hard in linear regression.

But can we Sigmoid fn \rightarrow (ReLU is better?)

$$y = \sigma(\underbrace{w^T x + b}_z)$$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

if z large $\rightarrow \sigma(z) \approx \frac{1}{1+0} = 1$.

if z large negative $\rightarrow \sigma(z) \approx \frac{1}{1+\infty} = 0$

Usually keep w and b separately

Functions you see: $x_0 = 1, x \in \mathbb{R}^{n_x+1}$

$$\hat{y} = \sigma(\underbrace{\theta^T x}_z)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \begin{array}{l} \rightarrow b \\ \} w. \end{array}$$

To implement a neural network, easy to keep w and b separate.

Ex 3: Logistic Regression Cost Function

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

learn $\{(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})\}$,

want $\hat{y}^{(i)} \approx y^{(i)}$

$\hat{y}^{(i)} \rightarrow$ data \rightarrow i th example
 $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$

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

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

Error (Loss) Function:

Common: Square Error: $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

Loss fun, to measure
how good our guess

but \downarrow not good for this b/c we
will get multiple optimum

is,

Write gradient descent.

For Logistic Regression we define a similar Loss Fun.

→ error for 1 example

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

Want to be as small as possible.

$$\text{if } y = 1: L(\hat{y}, y) = -\log \hat{y}$$

∴ $-\log \hat{y}$ to be as small as possible

⇒ $\log \hat{y}$ to be as large as possible

⇒ \hat{y} to be as large as possible (1).
 $\hat{y} \rightarrow$

$$\text{if } y = 0: L(\hat{y}, y) = -\log(1-\hat{y})$$

∴ $-\log(1-\hat{y})$ to be as small as possible

⇒ $\log(1-\hat{y})$ to be large.

⇒ $\underbrace{1-\hat{y}}_{\text{as big as possible}} \Rightarrow \hat{y} \text{ as small as possible}$
∴ $\hat{y} \rightarrow 0$

→ Avg over the entire training set.

Cost Function → How well you're doing over the
ENTIRE TRAINING SET

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n L(\hat{y}^{(i)}, y^{(i)})$$

\uparrow
Prediction

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

↳ minimize the cost function by varying tuning w, b .