# Deep learning
## Binary Classification Problems:
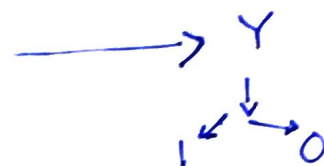
eg    Image  ———→  label: 0 or 1
      (Input)        (output label)



$X = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 233 \\ \vdots \\ 255 \\ 127 \end{bmatrix}$  ———→ Y
                                    $1 \swarrow \downarrow \searrow 0$

$64 \times 64$

$64 \times 64 \times 3 = 12288$          $12288$

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

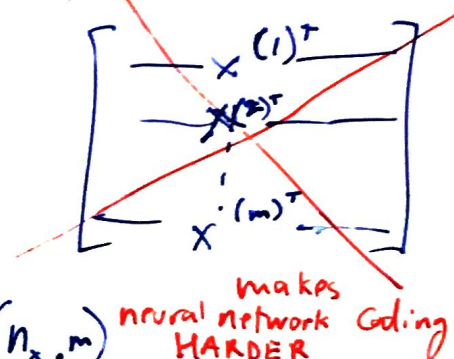  $m_{test}$ = # test examples = no. of test examples

$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix} \Big\updownarrow n_x$

$\xleftarrow{\hspace{2cm} m \hspace{2cm}}$

$\begin{bmatrix} \rule{1cm}{0.4pt} x^{(1)T} \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} x^{(2)T} \rule{1cm}{0.4pt} \\ \vdots \\ \rule{1cm}{0.4pt} x^{(m)T} \rule{1cm}{0.4pt} \end{bmatrix}$

makes neural network coding HARDER

X.shape = $(n_x, m)$

Similarly

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

## Logistic Regression:
  ——→ learning algorithm used when output labels (Y), in a supervised
  learning problem, are either 0 or 1   [ie for Bin. Class. Probs]
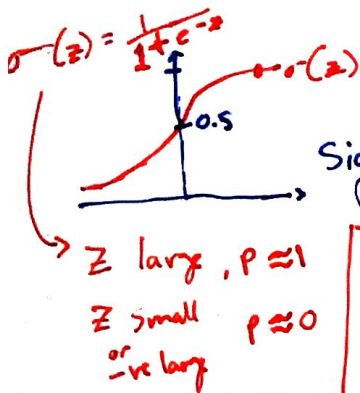
  Given x, want $\hat{y} = P(y=1 \mid x)$
          $x \in \mathbb{R}^{n_x}$
  Parameters:   $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

Output $\hat{y} = w^T x + b$  $\leftarrow$ — linear Regression!

∴ not to use here, it gives $y$ = actual value
we want $P(y=1|x)$

we'd want $\hat{y} = \sigma\left(\underbrace{w^T x + b}_{z}\right)$

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

0.5

$\sigma(z)$

Sigmoid fⁿ : gives values between 0 & 1

→ Z large, $p \approx 1$
Z small $p \approx 0$
or -ve large

Big Oof:
alternate notation
taught @ drexel ✱

$x_0 = 1$, $X \in \mathbb{R}^{n_x + 1}$

$\hat{Y} = \sigma(\theta^T x)$

✱ I used this notation in log. regression project.
learn this notation from now on.

$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix}$  } b → real no. or intercept

} w

$\underbrace{\qquad}$ we're not gonna use this notation
moving forward

QU: what's the parameters of logistic Regression?
Ans     W, an $n_x$ dimensional vector, & b, a real number

$\theta_{1 \to n}$  or  ,  $\theta_0$

→ x ——— x ——— x ——

✱ logistic regression can be viewed as a very small neural network

# logistic Regression Cost function:

→ to train the parameters $w$ & $b$ of logistic regression model, you need to define a cost function.

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

$$z^{(i)} = w^T x^{(i)} + b$$
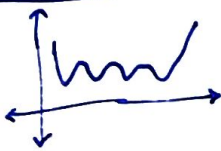
given $\{(x^{(1)}, y^{(1)}) \ldots (x^m, y^{(m)})\}$ we want $\hat{y}^{(i)} \approx y^{(i)}$

✳ **Loss (error function):**

1) You could do $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ → half square error

→ this basically measures how well our algo is doing

→ If you do this however, the optimization problem (later) becomes <u>non-convex</u> → end up with optimization problem with lots of <u>local optima</u>

↓

gradient descent might not find local optima



2) another way

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

rationale: if $y=1$ $L(\hat{y}, y) = -\log \hat{y}$ (← want as small as possible)
$\hat{y} = $ large → sigmoid cant be $>1$

if $y=0$ $L(\hat{y}, y) = -\log(1-y)$ want by $1-\hat{y}$ large want $\hat{y}$ small

✳ **Cost function**

→ loss f$^n$ was defined wrt single training example
→ Cost f$^n$ determines how well you're doing on entire training set

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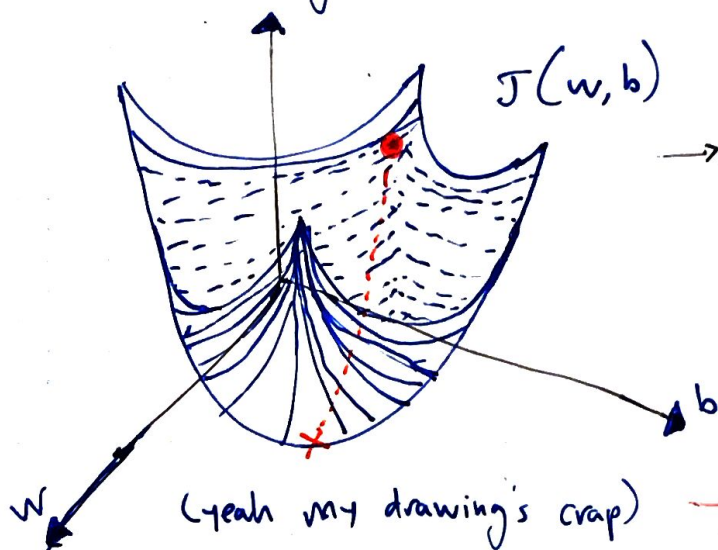
basically gives average of $L(\hat{y}, y)$

→ In training log. reg. model, we're going to find params $w$ & $b$ that minimize overall cost f$^n$ $J$.

# Gradient Descent:

Now let's actually try & use Gradient descent to learn params $w$ & $b$!

$J(w,b)$

→ want $w, b$ to minimize
$J(w,b)$

→ initialize $w, b$ to some initial val

→ for log reg, anything works,
→ random →0

(yeah my drawing's crap)

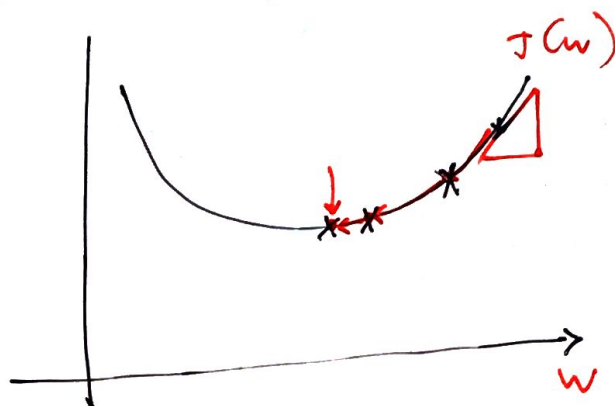→ gradient descent starts at that initial pt. & takes a step in steepest downhill direction

update, derivative

Repeat {
$$w : w - \alpha \frac{dJ(w)}{dw}$$

learning rate

}

$J(w)$

$w$

$$b := b - \alpha \frac{dJ(w,b)}{db}$$

$$w := w - \alpha \frac{(dJ(w,b))}{dw}$$

☺ also can write as $\frac{\partial J(w,b)}{\partial w}$

means the same thing mostly ~ior

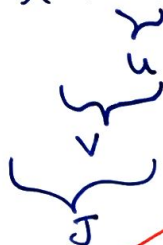⚡ use $\partial$ if 2 or more params → partial deriv
$d$ if one

# Side Quest: Computational Graphs    [Visualization trick/tool]

→ these come in handy when there is some distinguished or
   some special output variable,

   → # handy in visualizing components of a much complex
      equation that you're tying to optimize.
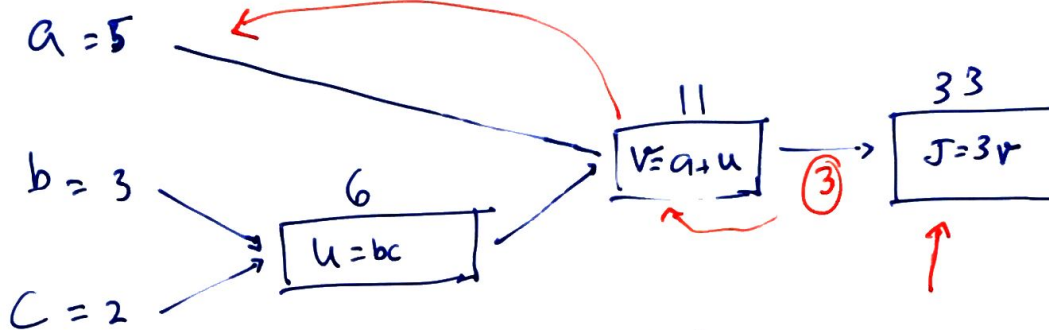
   → you can visualize each component to find bottlenecks!

$$J(a,b,c) = 3(a+bc) \qquad = 3(5+3\cdot2) = 33$$

assume {
  $\underbrace{\phantom{a}}_{u}$
  $\underbrace{\phantom{a}}_{v}$
  $\underbrace{\phantom{a}}_{J}$
}

Computational graph

$U = bc$
$v = a+u$
$J = 3v$

$a = 5$

$b = 3$

$C = 2$

| | 6 | | 11 | | 33 |
| $u=bc$ | | $v=a+u$ | ③ | $J=3v$ |

for computing derivatives we move R to L

---

$\dfrac{\partial J}{du} = ?$     $u = 6 \to 6.001$     $\Rightarrow \dfrac{dJ}{du} = 3$

$v = 11 \to 11.001$

$J = 33 \to 33.003$

i.e. $\dfrac{dJ}{dv} \cdot \dfrac{dv}{du}$
$\phantom{i.e.}\; \underset{3}{\underbrace{\phantom{dJ}}} \; \underset{1}{\underbrace{\phantom{dv}}}$

$J = 3v$

① $\dfrac{dJ}{dv} = ? \longrightarrow = 3$

② $\dfrac{dJ}{da} = a?$  → If we bump up a, whats effect on J?

$u = 5 \to 5.001$
$v = 11 \to 11.001$
$J = 33 \to 33.003$
$a \xrightarrow{\triangle} 0.001 \qquad J \xrightarrow{\triangle} 0.003$

$\dfrac{dJ}{da} = 3$

$\dfrac{dJ}{db} = \dfrac{dJ}{du} \cdot \dfrac{du}{db}$     $b = 3 \to 3.001$  ④
$u = b \cdot c = 6 \to 6.002$  $c=2$

$v = 11.002$
$J = 33.006$   $\dfrac{dJ}{db} = 6$

TIdr:  forward  L→R ⇒ compute cost $f^n$
       backward R→L ⇒ compute derivati

# Logistic Regression Gradient Descent

**for one observation:**     $\hat{y} = a = \sigma(z)$     $z = w^T x + b$

$X_1$

$W_1$

$X_2$ $\rightarrow$ $\boxed{z = w_1 X_1 + w_2 X_2 + b}$ $\rightarrow$ $\boxed{\hat{y} = a = \sigma(z)}$ $\leftarrow$ $\boxed{L(a,y)}$

$W_2$

$b$

In logistic Regression, we want to modify params $w \& b$ to reduce the <u>loss</u>

How you actually compute the loss on a single training eg,
we want to compute derivatives wrt this loss

1) go backward to compute derivative of loss wrt a
"da" $= \dfrac{dL(a,y)}{da}$

$$\boxed{= -\dfrac{y}{a} + \dfrac{1-y}{1-a}}$$

$$\boxed{dz = \dfrac{dL}{dz} = \dfrac{dL(a,y)}{dz}}$$
$\quad = a - y$
$\quad = \dfrac{dL}{da} \cdot \dfrac{da}{dz} \rightarrow a(1-a)$

$\dfrac{dL}{dw_1} = "dw_1" = X_1 \cdot dz$ ; $dw_2 = x_2 \cdot dz$ ; $db = dz$

$w_1 := w_1 - \alpha dw_1$   $w_2 := w_2 - \alpha dw_2$   $b := b - \alpha db$

**formula for derivative of loss wrt z?**    $a - y$

Note # this explaination is for a single observation (assuming 2 features $X_1, x_2$)

Logistic Regression on $m$ examples:

loss $\quad J(w,b) = \frac{1}{m} \sum\limits_{i=1}^{m} L(a^{(i)}, y)$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum\limits_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_i} L(a^{(i)}, y^{(i)})}_{\partial w_1^{(i)}} - (x^{(i)}, y^{(i)})$$

initiate

$J = 0 \quad dw_1 = 0, \, dw_2 = 0, \, db = 0$

$\quad$ for $i = 1$ to $m$

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

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

$\qquad J \mathrel{+}= -\left[ y^{(i)} \log a^{(i)} + ((1 - \hat{y}^{(i)}) \log (1 - a^{(i)})) \right]$

$\qquad dz^{(i)} = a^{(i)} - y^{(i)}$

$\qquad dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)}$

$\qquad dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \quad \left.\right\}$ assuming 2 features $n=2$

$\qquad db \mathrel{+}= dz^{(i)}$

$\quad J \mathrel{/}= m$

$\quad dw_1 \mathrel{/}= m, \, dw_2 \mathrel{/}= m \, ; \, db \mathrel{/}= m$

$\partial w_1 = \frac{\partial J}{\partial w_1}$

$w_1 := w_1 - \alpha \, dw_1$

$w_2 := w_2 - \alpha \, dw_2$

$b := b - \alpha \, db$

Vectorization