

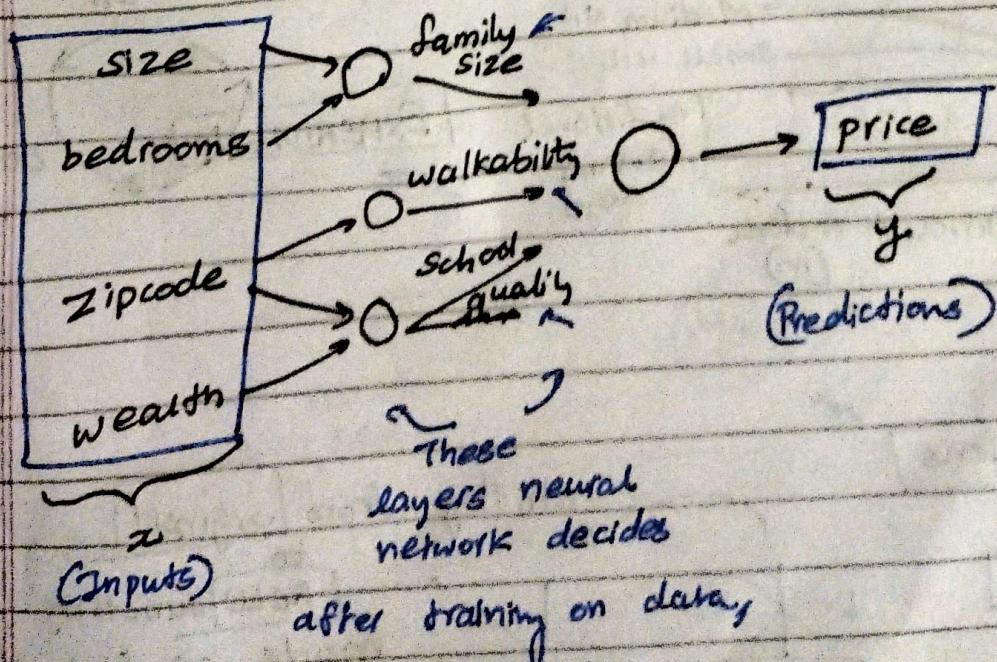
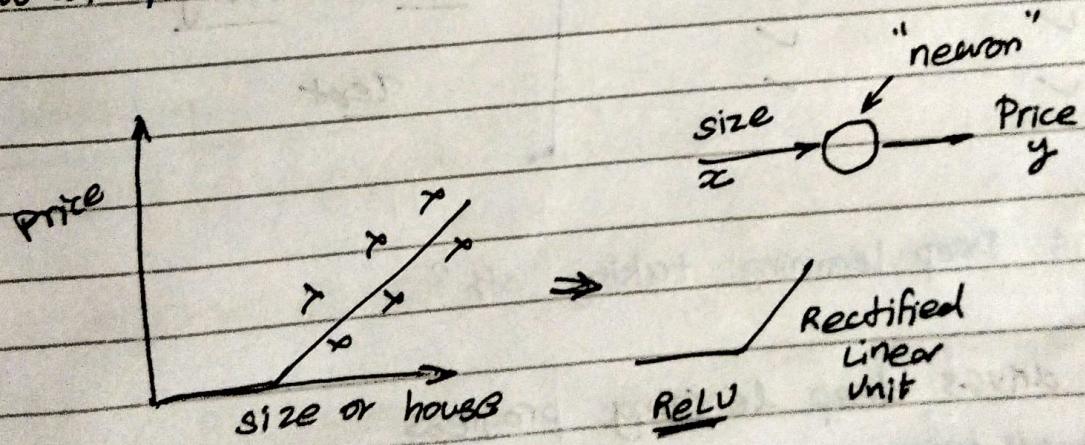
→ AI is new electricity (

→ as important as electricity revolution was

Courses

- Neural Networks & Deep Learning
- Improving deep learning Neural networks (Practical)
- Structuring your Machine learning project
- Convolutional Neural Networks.
- Natural language processing: Building sequence models

* what is a Neural Network?



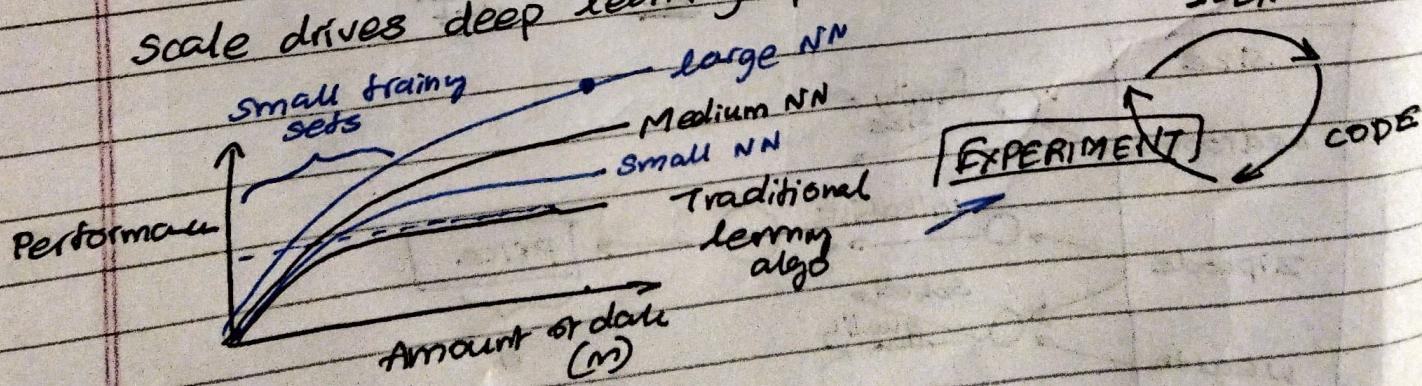
* Supervised learning with neural networks.

	INPUT (x)	OUTPUT (y)	
Standard NN	Home features Ad user info	price click on ad?	Real Estate online Advertising
CNN	Image	object (1...1000)	Photo Tagging
RNN	Audio English	Text transcript chinese	Speech Recognition Machine Translation
Custom/ Hybrid	Image, Radar Info	Position of other cars	Autonomous driving

	Structured Data	Unstructured Data
size	# bedrooms	Price
fixed	✓	✓
	✓	✓
	✓	✓
	✓	✓

* Why is deep learning taking off?

Scale drives deep learning progress



- Data

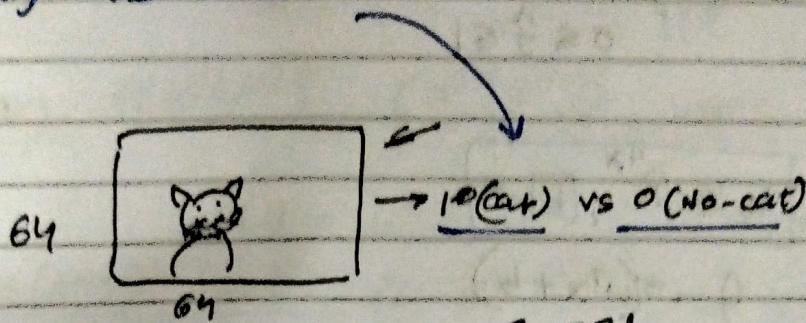
- computations

shifting from sigmoid

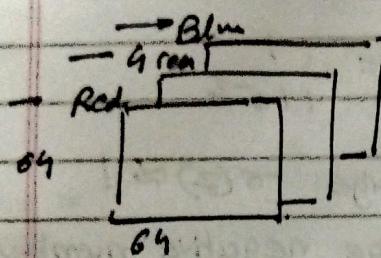
to ReLU function
made gradient descent work much faster

* Logistic Regression as a Neural Network.

* Binary classification



$\rightarrow 1 \text{ (cat)} \text{ vs } 0 \text{ (no-cat)}$



$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 137 \end{bmatrix} \left\{ \begin{array}{l} \text{Red} \\ \text{Green} \\ \text{Blue} \end{array} \right. \quad 64 \times 64 \times 3 = 12288$$

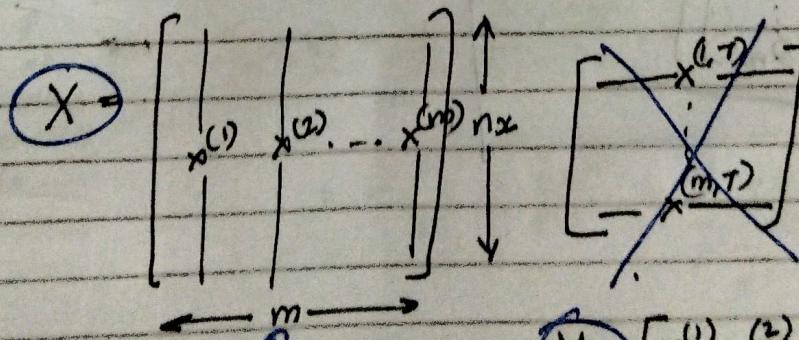
$$n = n_x = 12288 \quad x \rightarrow y$$

Notation

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

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

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



much easier

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

$$X.\text{shape} = (n_x, m)$$

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

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

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

Logistic Regression

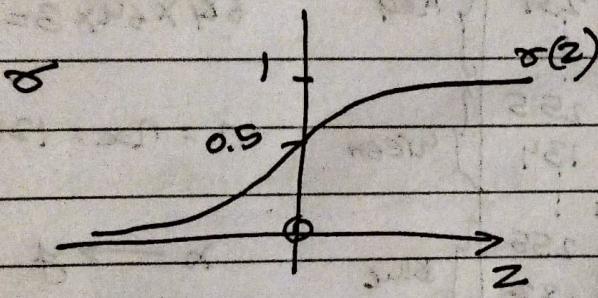
Given x , want $\hat{y} = P(y=1|x)$

$$0 \leq \hat{y} \leq 1$$

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

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

output: $\hat{y} = \sigma(w^T x + b)$



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

If z large $\sigma(z) \approx 1$

If z large negative number

$$\sigma(z) = \frac{1}{1+\infty} \approx 0$$

Different convention

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_a} \end{bmatrix} \left\{ \begin{array}{l} b \leftarrow \\ w \leftarrow \end{array} \right.$$

Not
using
these

Logistic Regression cost function

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

you should have got

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

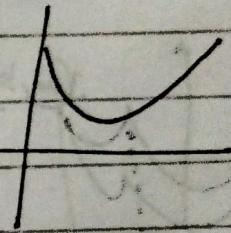
that you got

$$z^{(i)} = w^T x^{(i)} + b^{(i)} \quad i^{\text{th}} \text{ example}$$

→ loss (error) function: (single training example)

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

↑ ↑
output end up with
optimisation problem



so this is used

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

↓ close

$$\begin{aligned} \text{if } y=1 & \quad L(\hat{y}, y) = -\log \hat{y} \leftarrow \text{want } \log \hat{y} \text{ large, want } \hat{y} \text{ large} \\ \text{if } y=0 & \quad L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{want } \log(1-\hat{y}) \text{ large} \end{aligned}$$

$L(\hat{y}, y) \rightarrow$ should be as small as possible

⇒ cost function: (cost of parameters)

$$\text{Average of } \overset{n}{\underset{i=1}{\sum}} \text{ loss functions over a entire training set}$$

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

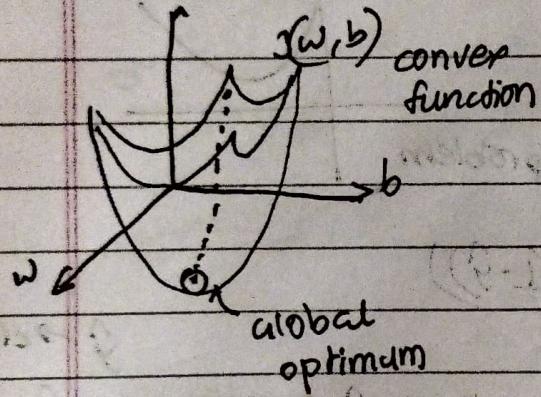
$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$

* Gradient Descent

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

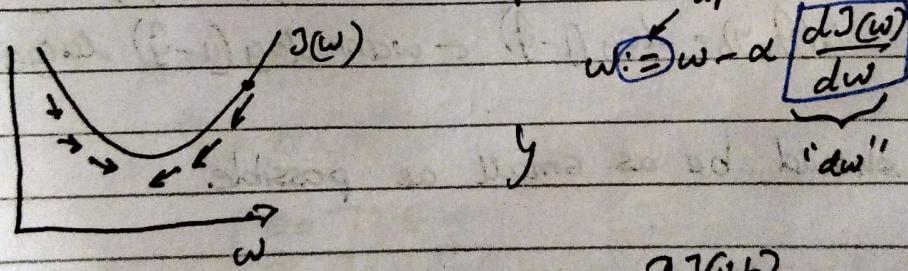
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

want to find w, b such that minimize $J(w, b)$



Formula

Repeat \leftarrow updates



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

"dw"

in 3ds $\left\{ \begin{array}{l} w := w - \alpha \frac{\partial J(w, b)}{\partial w} \\ b := b - \alpha \frac{\partial J(w, b)}{\partial b} \end{array} \right.$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

y

Derivatives

Kya Dekhu

BOPM

classmate

Date _____

Page _____

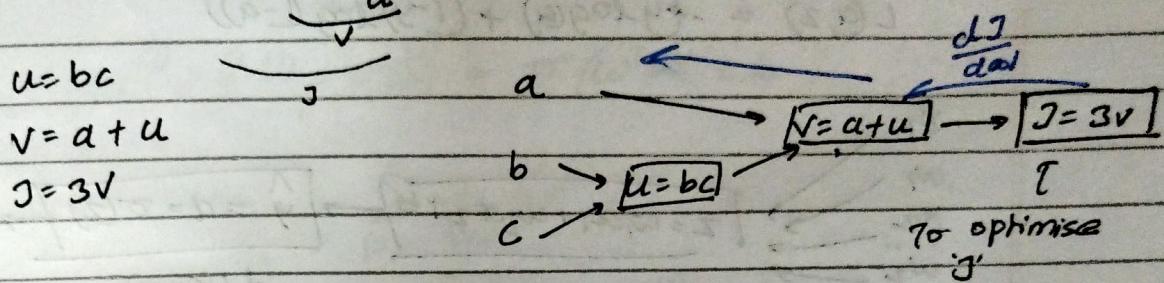
* Computation Graph

$$J(a, b, c) = 3 \underbrace{(a + bc)}_u$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



* Derivatives with a computation graph

$$J = 3v$$

$$\frac{dJ}{dv} = 3 \quad \therefore "dv" = 3$$

$$\begin{aligned} \frac{dJ}{da} &= \frac{dv}{da} \frac{dJ}{dv} \quad \therefore "da" = 3 \\ &= 3 \times 1 = 3 \end{aligned}$$

Final output variable J ,

many variables still will

$$\frac{d \text{Final output var}}{d \text{var}} \rightarrow "d \text{var}"$$

similarly $du = 3$

$$db = 6$$

$$\frac{dJ}{db} = \frac{dJ}{du} \frac{du}{db} = 6$$

3 2

Logistic Regression Gradient Descent

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

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

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

$$\begin{aligned}
 & x_1 \\
 & w_1 \\
 & x_2 \\
 & w_2 \\
 & b
 \end{aligned}
 \rightarrow
 \boxed{z = w_1 x_1 + w_2 x_2 + b} \rightarrow
 \boxed{\hat{y} = a = \sigma(z)} \rightarrow L(a, y)$$

$$\begin{aligned}
 dz &= \frac{dL}{dz} = \frac{dL(a, y)}{dz} & "da" &= \frac{dL(a, y)}{da} \\
 &= a - y & &= -\frac{y}{a} + \frac{1-y}{1-a} \\
 &= \underbrace{\frac{da}{da}}_{a(1-a)} \underbrace{\frac{da}{dz}}_{dL/dz}
 \end{aligned}$$

$$\frac{dL}{dw_1} = "dw_1" = x_1 dz \quad dw_2 = x_2 dz \quad db = dz$$

$$w_1 = w_1 - \alpha dw_1$$

$$w_2 = w_2 - \alpha dw_2$$

$$b = b - \alpha db$$

Derivation

$$\begin{aligned}
 & \cancel{\left(\frac{1}{1+e^{-z}}\right)^2} \\
 & \cancel{\left(\frac{e^{-z}}{1+e^{-z}}\right)^2} \quad \left(\frac{2}{1+e^{-z}}\right) \\
 & -\cancel{\frac{(1+e^{-z})+1}{(1+e^{-z})^2}} = -\frac{1}{1+e^{-z}} + \frac{1}{(1+e^{-z})^2} \\
 & \frac{1+e^{-z}+2(e^{-z})}{(1+e^{-z})^2} = -a + a^2 \cancel{a(a-1)}
 \end{aligned}$$

$$\begin{aligned}
 & \frac{1}{1+e^{-z}} + \frac{2e^{-z}}{(1+e^{-z})^2} \\
 & = \frac{1}{1+e^{-z}} + \frac{2e^{-z}}{(1+e^{-z})^2}
 \end{aligned}$$

* Logistic Regression on m examples

$$J(\underline{w}, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

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

$$\frac{\partial}{\partial w_j} J(\underline{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_j} L(a^{(i)}, y^{(i)})}_{dL^{(i)} - (a^{(i)}, y^{(i)})}$$

Initialise $J=0$; $d\underline{w}_1=0$; $d\underline{w}_2=0$; $db=0$

For $i=1$ to m

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

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

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

$$\begin{aligned} d\underline{w}_1 &+= x_1^{(i)} dz^{(i)} \\ d\underline{w}_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned} \quad \uparrow n=2$$

$$j/m$$

$$d\underline{w}_1/m; \quad d\underline{w}_2/m; \quad db/m$$

1

$$w_1 := w_1 - \alpha d\underline{w}_1$$

$$w_2 := w_2 - \alpha d\underline{w}_2$$

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

$$d\underline{w}_1 = \frac{dJ}{d\underline{w}_1}$$

Vectorization

Art of getting rid of explicit for loop in code

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

$$w = \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} \quad x = \begin{bmatrix} | \\ | \\ | \\ | \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Non vectorized:

$$z = 0$$

for i in range(n-x);

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

$$z += b$$

than this

vectorized:

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

Takes much less time

GPU } SIMD — single instruction
CPU } multiple data

→ whenever possible, avoid explicit for loops

$$u = Av$$

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

$$u_i = \sum_j A_{ij} v_j$$

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

for i - -

for j - -

$$u[j] += A[i][j] * v[j]$$

$$u = np.dot(A, v)$$

$$np.log(v)$$

$$np.abs(v)$$

$$np.maximum(v, 0)$$

$$v ** 2$$

$$\frac{1}{v}$$

Instead of $dw_1 = 0$, $dw_2 = 0$

$$dw = \text{np.zeros}((n_x, 1))$$

one for loop removed

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array} \right\} \boxed{dw += x^{(i)} dz^{(i)}} \text{ vectors}$$

$$\left. \begin{array}{l} dw_1 = \frac{dw_1}{m} \\ dw_2 = \frac{dw_2}{m} \end{array} \right\} \boxed{dw = m}$$

Vectorizing Logistic Regression

$$\begin{aligned} \rightarrow z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

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

$$R^{(n_x, m)} \quad [z^{(1)} \ z^{(2)} \dots z^{(m)}] = w^T x + [b \ b \ b \dots b]_{1 \times m}$$

$$\begin{aligned} w^T x &= [w^T x^{(1)} + w^T x_2 + w^T x_3 + \dots] [b \ b \ b] \\ &= [w^T x^{(1)} + b + w^T x^{(2)} + b + w^T x^{(3)} + b + \dots] \end{aligned}$$

$$\Rightarrow \boxed{z = np.dot(w.T, x) + b} \rightarrow \begin{array}{l} \text{Broadcasting takes} \\ \text{place \&} \\ \text{this becomes} \end{array}$$

$$[b \ b \ b \dots b]$$

$$a = [a^{(1)} \ a^{(2)} \dots a^{(m)}] = \sigma(z)$$

$$dz^{(i)} = a^{(i)} - y^{(i)} \quad dz^{(0)} = a^{(0)} - y^{(0)}$$

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

$$A = [a^{(1)} \dots \ a^{(m)}] \quad y = [y^{(1)} \dots \ y^{(m)}]$$

$$dz = A - y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots]$$

$$dw = 0$$

$$db = 0$$

$$\begin{bmatrix} dw^{(1)} = x^{(1)} dz^{(1)} \\ dw^{(2)} = x^{(2)} dz^{(2)} \\ \vdots \\ dw^{(m)} = x^{(m)} dz^{(m)} \end{bmatrix} \quad \begin{bmatrix} db^{(1)} = dz^{(1)} \\ db^{(2)} = dz^{(2)} \\ \vdots \\ db^{(m)} = dz^{(m)} \end{bmatrix}$$

$$dp = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

$$= \frac{1}{m} \text{np.sum}(dz)$$

$$dw = \frac{1}{m} X dz^T$$

$$= \frac{1}{m} \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]$$

complement by using vectors

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

$$= \text{np.dot}(w^T, x) + b$$

$$A = \sigma(z)$$

$$dz = A - y$$

$$dw = \frac{1}{m} X dz^T$$

$$db = \frac{1}{m} \text{np.sum}(dz)$$

$$w := w - \alpha dw$$

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

\rightarrow For iteration & whole
for loop req,

Broadcasting in Python

	Apples	Beet	Eggs	Potato	
Carb	56	0.0	4.4	68.0	
Protein	1.2	104.0	52.0	8.0	
Fat	1.8	135.0	99.0	0.9	

$= A$
 $(3,4)$

calculate % of calories from C, P, F:
can you do this without for loop?

$cal = A \cdot \text{sum}(axis=0) \leftarrow \text{sum vertically}$
 print(cal)
 $[59, 239, 155.4, 76.9]$

② percentage = $100 * A / \text{cal} \cdot \text{reshape}(1,4)$
 $1(3,4) / (4,1)$ used to reshape matrix

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \frac{1}{100} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

↓

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \downarrow \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

General Principle

$$\begin{array}{ccc}
 (m,n) & + & (1,n) \rightsquigarrow (m,n) \\
 \text{matrix} & \times & \\
 & / & \\
 & & (m,1) \rightsquigarrow (m,n)
 \end{array}$$

$$(m,1) + R$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

Matlab/Octave: bsxfun ← broadcast

Python/Numpy

Tips & tricks to eliminate bugs

Specify rows and columns always as
this may create a ambiguity b/w array & matrix

a = np.random.randn(5)
 a.shape = (5,) } Beware
 "rank 1 array" } Don't use

a = np.random.randn(5, 1) → column vector
 a = np.random.randn(1, 5) → Row vector

assert (a.shape == (5, 1)) ←

a = a.reshape((5, 1)) ← To convert as per need

Jupyter/Python Notebooks

→ Shift + Enter to Execute program

- Restart kernel if get any errors
- submit assignment

Explanation of Logistic Regression cost function

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

Interpret

$$\begin{cases} \hat{y} = P(y=1/x) \\ \text{if } y=1 : P(y/x) = \hat{y} \\ y=0 : P(y/x) = 1-\hat{y} \end{cases}$$

$$P(y/x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\log P(\hat{y}/x) = \log y \log \hat{y} + (1-y) \log (1-\hat{y})$$

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

cost on m examples

$$\log P(\text{labels in training set}) = \sum_{i=1}^m -L(\hat{y}^{(i)}, y^{(i)})$$

$$\log P(\dots) = \underbrace{\sum_{i=1}^m \log P(\hat{y}^{(i)}/x^{(i)})}_{-L(\hat{y}^{(i)}, y^{(i)})}$$

maximum likelihood estimation

$$\text{cost } J(w, b) = -\sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

Cost

$$\rightarrow \left[\begin{array}{c|c} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ \hline 1 & 1 \end{array} \right] \quad \left\{ \begin{array}{l} \hat{y} \\ y \end{array} \right\}$$

0.5 0.1