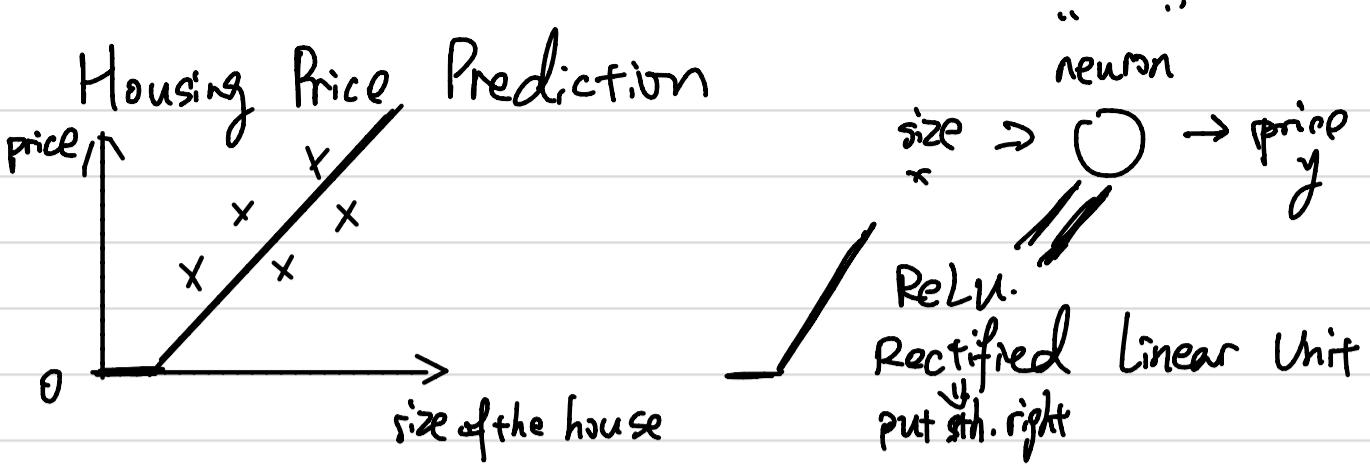
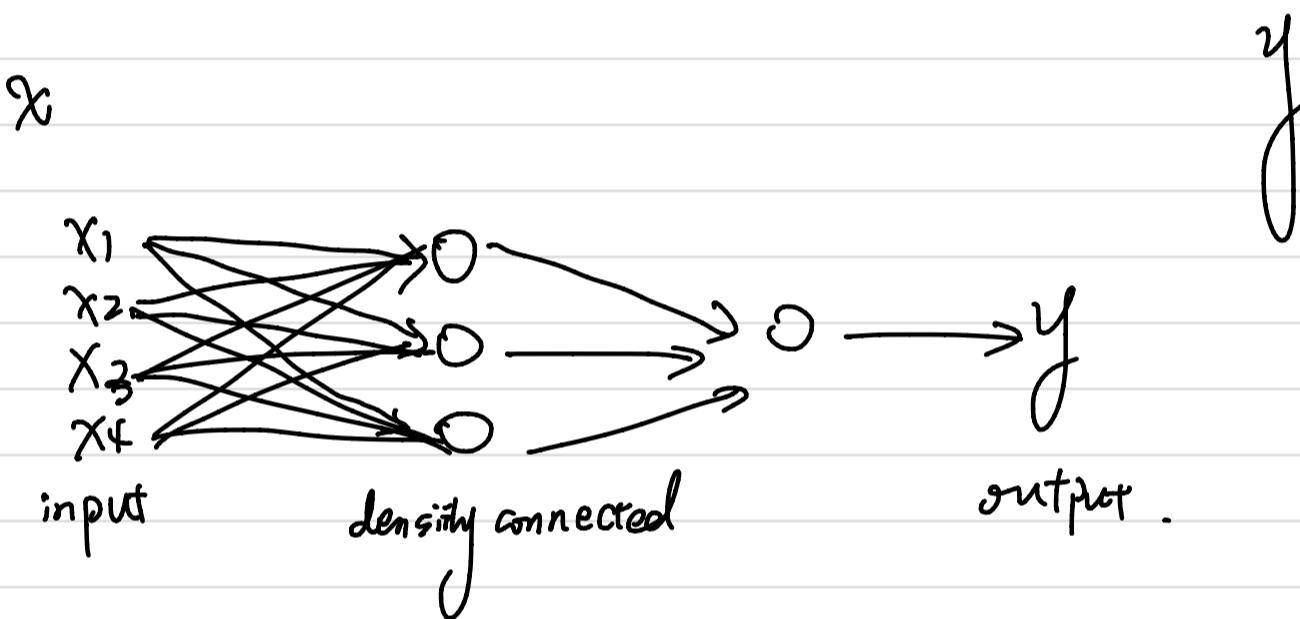
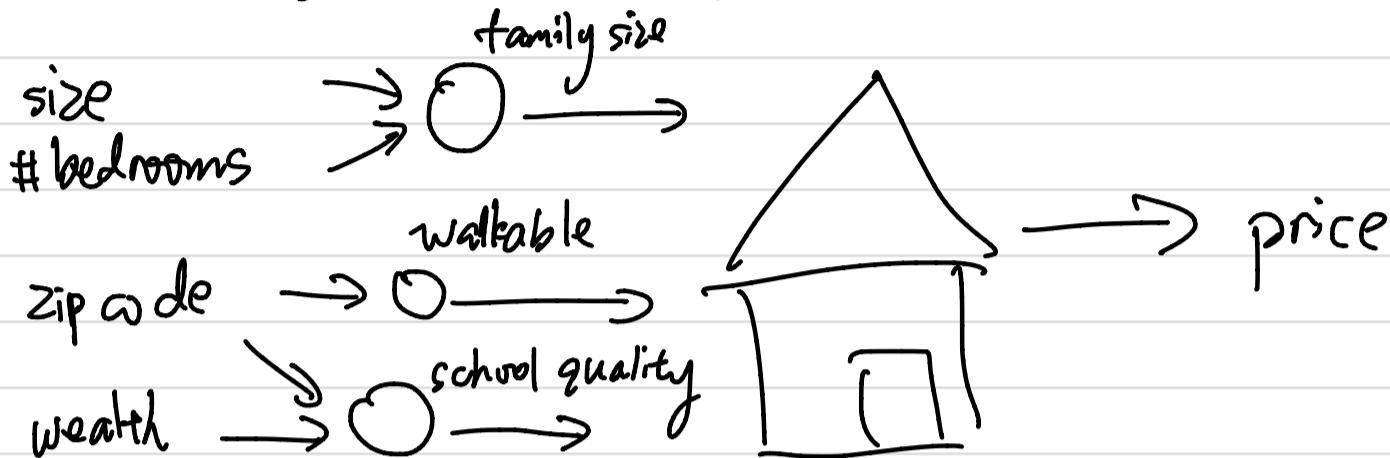


Andrew Ng.

Deep learning class on Coursera



Here neuron just this linear regression. we have thousands of neurons.



Supervised Learning

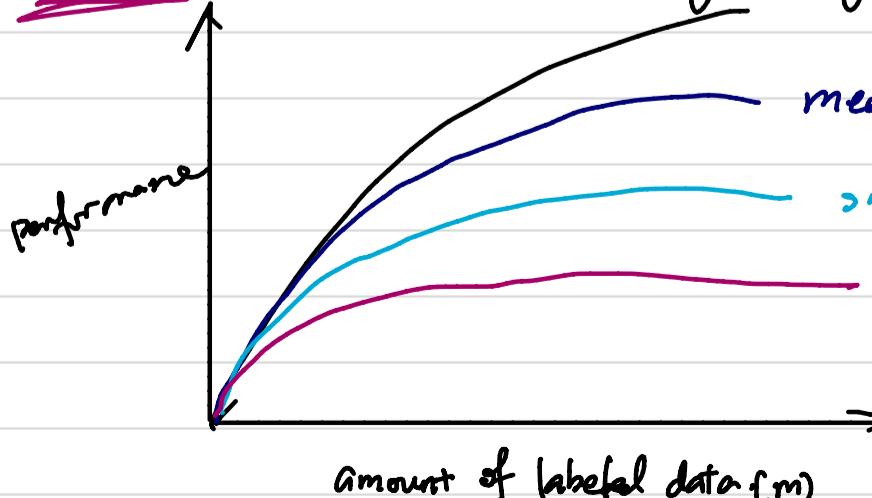
convolutional neural network for images.

Machine learning.

structured Data
feature has defined meanings.
i.e. number, price, size

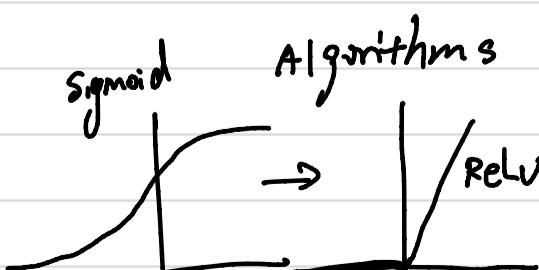
unstructured data
audio, image, text

scale drives deep learning progress



Data

computation



Logistic Regression

• Binary classification 1 cat vs 0 not cat

* Training data is implemented to build up a model, while a test data set is to validate the model built.

Logistic Regression

Given x . y : be a cat.

$$\hat{y} = p(y=1|x)$$

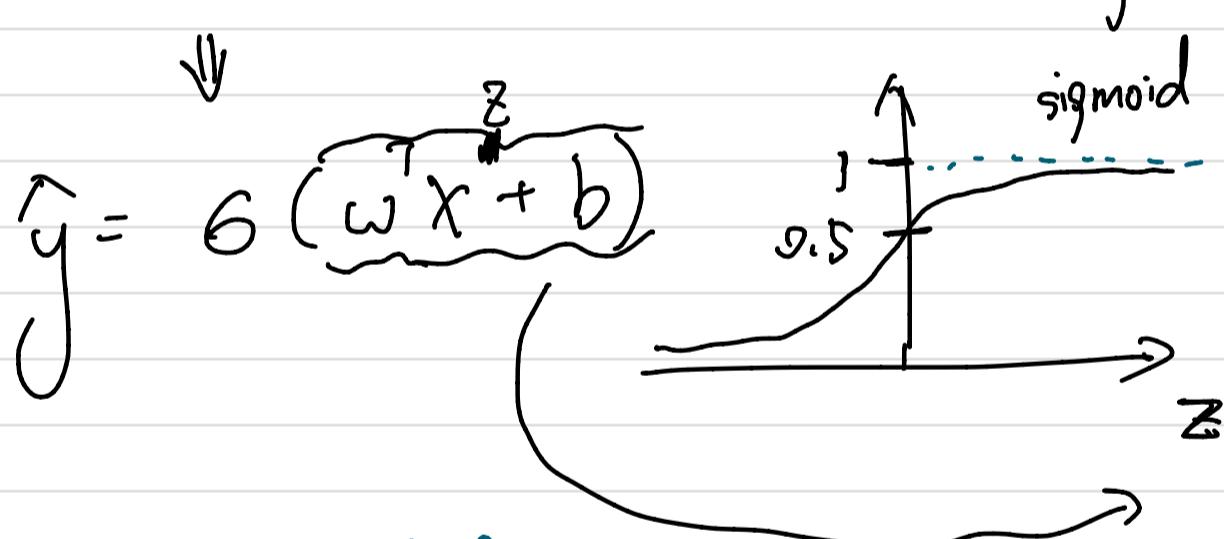
given x . what's the prob of y is cat

$x \in \mathbb{R}$

Parameter : $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$. $0 \leq \hat{y} \leq 1$

Output

$\hat{y} = w^T x + b$ this value can be negative or very big, which contradicts $0 \leq \hat{y} \leq 1$



$z = w^T x + b$ 是一个
酱油. $\hat{y} = g(z)$ 才
是我们要估计的 y 值.

Formula for sigmoid function:

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

If z large, $g(z) \approx \frac{1}{1+0} = 1$

If z small, $g(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Big number}} \approx 0$ value from the algorithm

Cost function

training data : Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want

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

the real training y value

Loss function / error function: single training example.

measure how good our output \hat{y} is when the true label is y .

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2 \text{ fit. 因为真实情况是 } \approx \text{ not } V.$$

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

want to as small as possible.

If $y=1$. $L(\hat{y}, y) = -\log \hat{y}$ ^{small} → want $\log \hat{y}$ small, want \hat{y} small

If $y=0$. $L(\hat{y}, y) = -\log(1-\hat{y})$ ^{small} → want $\log(1-\hat{y})$ small. want \hat{y} ^{large}

Cost function

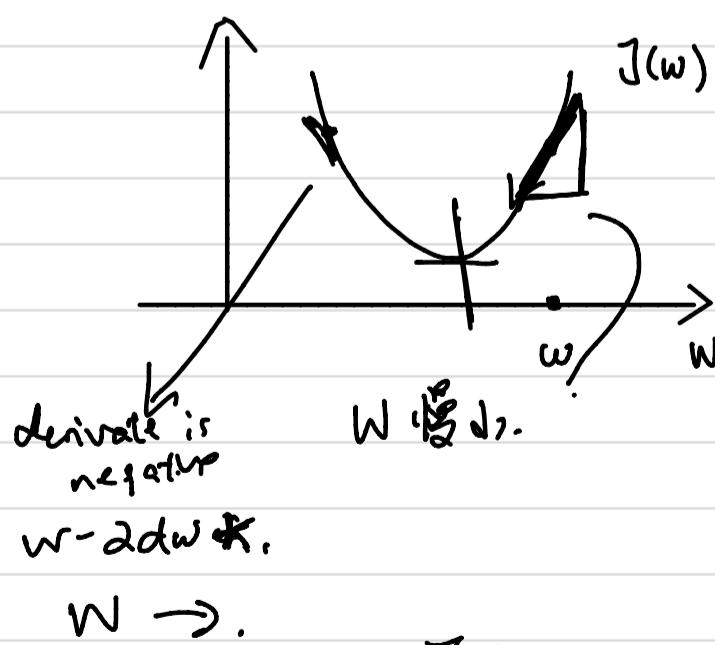
measure how good for entire ^{training} data set.

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

we want to find
the w, b to minimize the
error for the whole training data

for convex

Gradient Descent



ignore b for now

Gradient
Repeat

$w :=$ update to

learning rate, controls how big each step we take

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

code: "dw"

derivative

what we write in code

symbol.

$\frac{\partial J(w, b)}{\partial w}$ J is a function of 2 variables

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

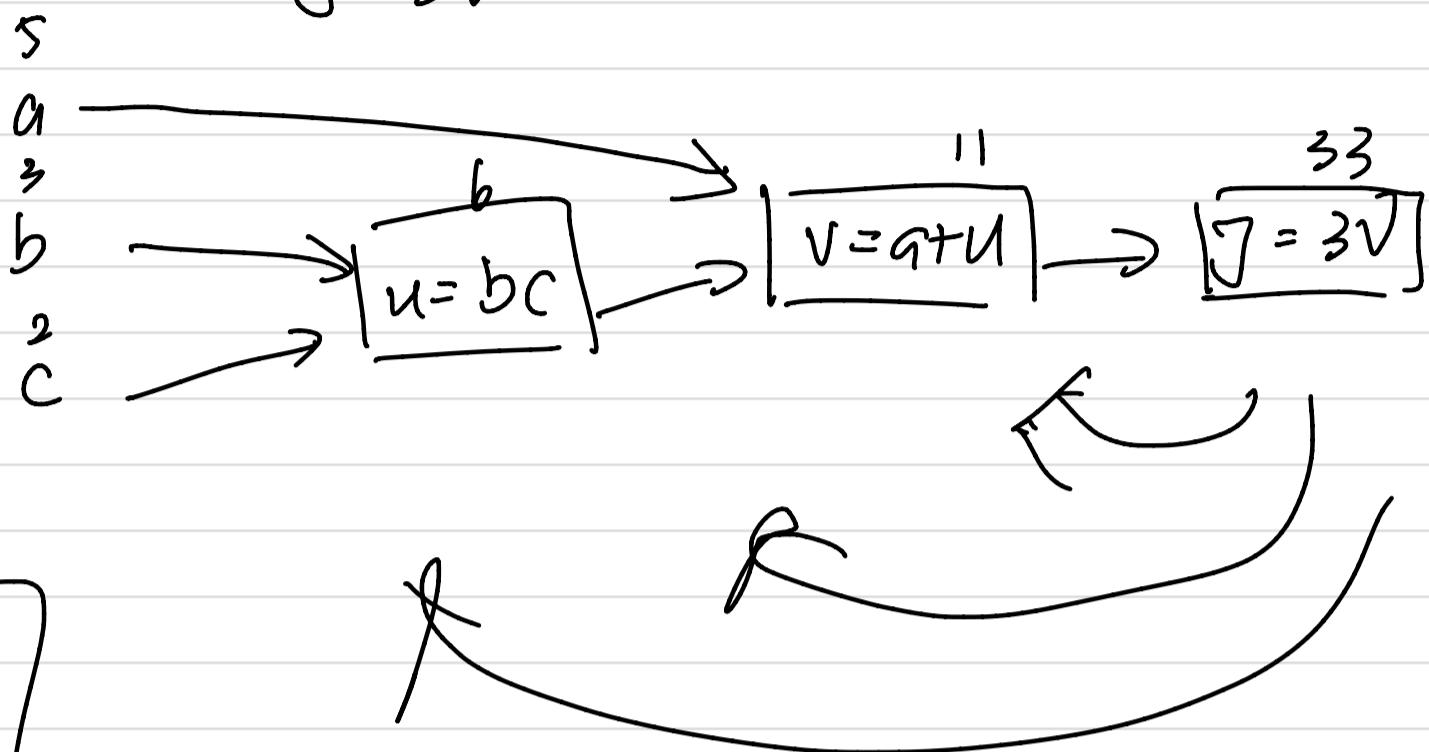
$$b := b - \alpha \frac{dJ(w, b)}{db}$$

$J(w, b)$

Computation graph.

$$J(a, b, c) = 3(a + bc)$$

$$\begin{aligned} U &= bc \\ V &= a + U \\ J &= 3V \end{aligned}$$



$$\boxed{\frac{dJ}{dV} = 3}$$

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

$$\boxed{\frac{dV}{da} = 1}$$

$$\frac{dJ}{da} = \frac{dJ}{dV} \cdot \frac{dV}{da}$$

change $a \rightarrow$ change $V \Rightarrow J$

$$a = 5 \rightarrow \underline{5.00}$$

$$V = 11 \rightarrow \underline{11.00}$$

$$J = 33 \rightarrow \underline{33.00}$$

\Rightarrow chain rule

$$\frac{dJ}{da} = \frac{dJ}{dV} \cdot \frac{dV}{da}$$

final output variable = coding.

(here is J)

$\frac{d \text{Final OutVar}}{d \text{var}}$ in code "dvar"

i.e.

$\frac{dJ}{da}$ in code "da"

$$\frac{dJ}{du} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{du}$$

3 1

$u = 6 \rightarrow 6.001$
 $v = 11 \rightarrow 11.001$
 $J = 33 \rightarrow 33.003$

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

3 2

$b = 3 \rightarrow 3.001$
 $u = b \cdot c \rightarrow 6.002$ ($c = 2$)

(in code "db")

$$\frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = 9.$$

3 3

Logistic Regression Recap.

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

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

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

loss function

example:

only 2 features x_1, x_2 :
 w_1, b w_2, b

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

~~未詳説、値の範囲~~

$$\frac{dz}{dZ} = \frac{dL}{dZ} = \frac{dL(a, y)}{dZ}$$

$= \frac{a - y}{a(1-a)}$

$$\therefore da'' = \frac{da}{dZ}$$

$= \frac{-y}{a} + \frac{1-y}{1-a}$

$$= \frac{dL}{da} \cdot \frac{da}{dZ} \rightarrow a(1-a) \left(\frac{-y}{a} + \frac{1-y}{1-a} \right) \cdot a(1-a)$$

最後の項は

$$dw_1 = x_1 \cdot dz$$

$$dw_2 = x_2 \cdot dz$$

$$db = dz$$

$$\frac{dL}{dw_1} = dw_1 = x_1 \cdot dz$$

$$dw_2 = x_2 \cdot dz$$

$$db = dz$$

update

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

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

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

Logistic on m examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) \quad \text{where } a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})$$

example:

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

For $i=1$ to m ①

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

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

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

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

$$dw_1 += x_1^{(i)} dz^{(i)} \quad n=2 \quad \text{②}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J / = m \quad \text{f. average} \quad dw_1 / = m$$

$$dw_2 / = m$$

$$dw_1 = \frac{\partial J}{\partial w_1}$$

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

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

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

(explicit)

✗ for long needs more time to execute
(long for loop)

Your code needs to be efficient \Rightarrow vector

$$dw_1^{(i)}, dw_2^{(i)}$$

$$dw_1^{(i)} - (x^{(i)}, y^{(i)})$$

~~for attribute~~ ~~M~~
~~for data~~ ~~for i~~
~~for j~~ ~~for k~~

2 for loops

(i) data ~~for~~
j : 1 ~ m

n 是 feature 数。

Vectorization

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

$$w = \begin{bmatrix} : \\ : \end{bmatrix}$$

$$x = \begin{bmatrix} : \\ : \end{bmatrix}$$

$$\begin{aligned} x &\in \mathbb{R}^n \\ x &\in \mathbb{R}^n \end{aligned}$$

Non-vectorize

```
z = 0
for i in range(n - x):
    z += w[i] * x[i]
z += b
```

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

Code:

```
import numpy as np
a = np.array([1, 2, 3, 4])
print(a)
```

```
import time
a = np.random.rand(1000000)
b = np.random.rand(1000)
tic = time.time()
c = np.dot(a, b)
toc = time.time()
print("Vectorized version: " + str((toc - tic) * 1000) + "ms")
```

Avoid explicit for loop

$$V = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \quad U = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

non-vectorization

```
U = np.zeros((n, 1))
for i in range(n):
    u[i] = math.exp(v[i])
```

vectorization

```
import numpy as np
u = np.exp(v)
* np.log(v)
* np.abs(v)
np.maximum(v, 0)
```

Broadcasting

$m \uparrow$ data $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ (attributes)

$$z^{(1)} = w^T x^{(1)} + b, \quad z^{(2)} = w^T x^{(2)} + b, \quad z^{(3)} = w^T x^{(3)} + b.$$

$$y^{(1)} = f(z^{(1)}) \quad a^{(2)} = f(z^{(2)}) \quad a^{(3)} = f(z^{(3)})$$

$$Z = [z^{(1)}, z^{(2)}, \dots, z^{(m)}] = W^T X + [b, b, \dots, b]_{1 \times m} = [W^T X^{(1)} + b, \dots, W^T X^{(m)} + b]$$

$Z = np.dot(W.T, X) + b$ 只需写 b. 自动生成 vector 由 numpy. "Broadcasting"

$$A = [a^{(1)}, a^{(2)}, \dots, a^{(m)}]$$

Review:

$$\therefore dz = a - y \\ dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \dots \dots$$

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

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

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

From before :

$dw = 0$
$dw^+ = X^{(1)} dz^{(1)}$
$dw^+ = X^{(2)} dz^{(2)}$
\vdots
$dw/m = m$

$db = 0$
$db^+ = dz^{(1)}$
$db^+ = dz^{(2)}$
\vdots
db/m

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

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

$$dw = \frac{1}{m} X dz^T \\ = \frac{1}{m} \left[\begin{array}{c|c|c|c} 1 & & & \\ \hline x^{(1)} & \cdots & x^{(m)} & \\ \hline \vdots & & \vdots & \\ \hline \end{array} \right] \left[\begin{array}{c} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{array} \right]_{m \times 1}$$

$$= \frac{1}{m} \left[X^{(1)} dz^{(1)} + \dots + X^{(m)} dz^{(m)} \right].$$

$$Z = W^T X + b$$

$$= np.dot(W.T, X) + b$$

$$A = g(Z)$$

$$dz = A - Y$$

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

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

$$w = w - \alpha \overbrace{dw}^{\text{learning rate}}$$

$$b = b - \alpha db$$

$A = np.array([[-\dots],$
 $[J, J]])$
 $cal = A.sum(axis=0)$. -3|3相加。

Broadcasting example

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix} + \underbrace{100}_{\text{in python will be}} \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \Rightarrow \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \\ 105 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}_{(1,n)} \Rightarrow \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

python: copy to turn into (m, n)

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{(m,1)} \stackrel{100}{\overbrace{\dots}}_{200} \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

$$\begin{array}{ccc} (m, n) & \xrightarrow{\text{matrix}} & (1, n) \\ & / & \\ & & (m, 1) \end{array} \quad \begin{array}{c} \xrightarrow{\text{copy}^{*m}} \\ \rightsquigarrow (m, n) \\ \xrightarrow{\text{copy}^{*n}} (m, n) \end{array}$$

$$(m, 1) + \underbrace{100}_{\text{copy}^{*m}} = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

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

More About Cost Function

Here we only talk about logistic regression! So that's why $y=1$ or $y=0$ 哪两种可能!

$$\hat{y} = g(\omega^T X + b) \text{ where } g(z) = \frac{1}{1+e^{-z}}$$

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

$$\left. \begin{array}{l} \text{If } y=1 : P(y|x) = \hat{y} \\ \text{If } y=0 : P(y|x) = 1 - \hat{y} \end{array} \right\} P(y|x).$$

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

可看作 Bernoulli (非成功即失败)

$$\text{If } y=1, P(y|x) = \hat{y} \cdot (1-\hat{y})^0 = \hat{y}$$

$$\text{If } y=0, P(y|x) = \hat{y}^0 (1-\hat{y})^1 = 1 \times (1-\hat{y}) = 1-\hat{y}$$

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

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

to find parameter to make this as big as possible

maximize likelihood.

$$\text{LP (labels in training set)} = \prod_{i=1}^m P(y^{(i)} | x^{(i)})$$

+ log

\Rightarrow

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

$$= -\sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

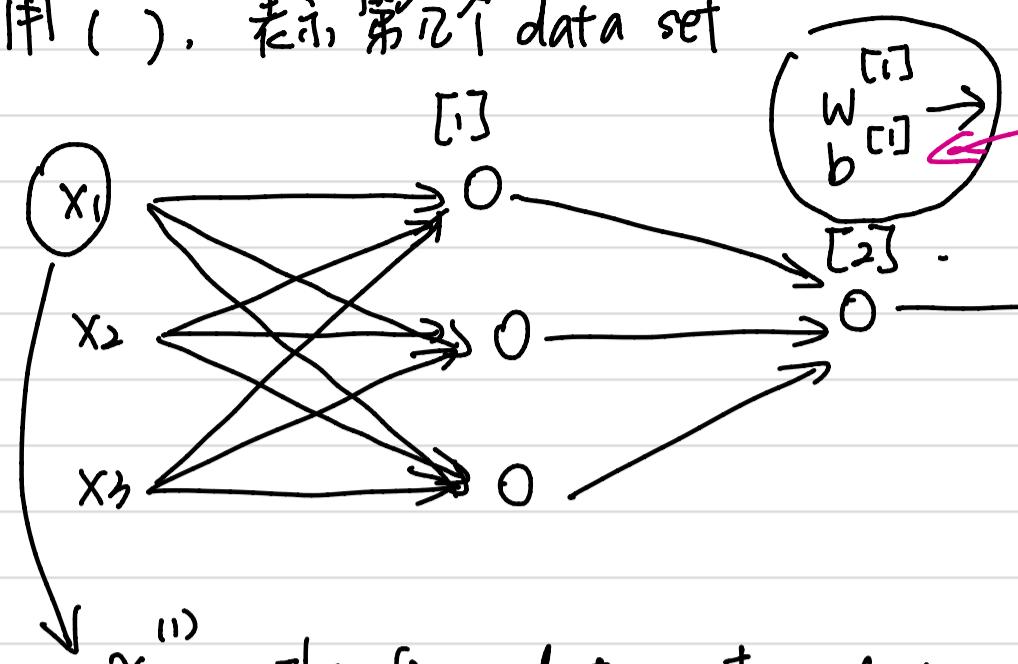
去完负号，要越小越好。

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

average, scaling

(minimize) 去负号。

用 []. 表示 layer
用 (). 表示 第 i 个 data set



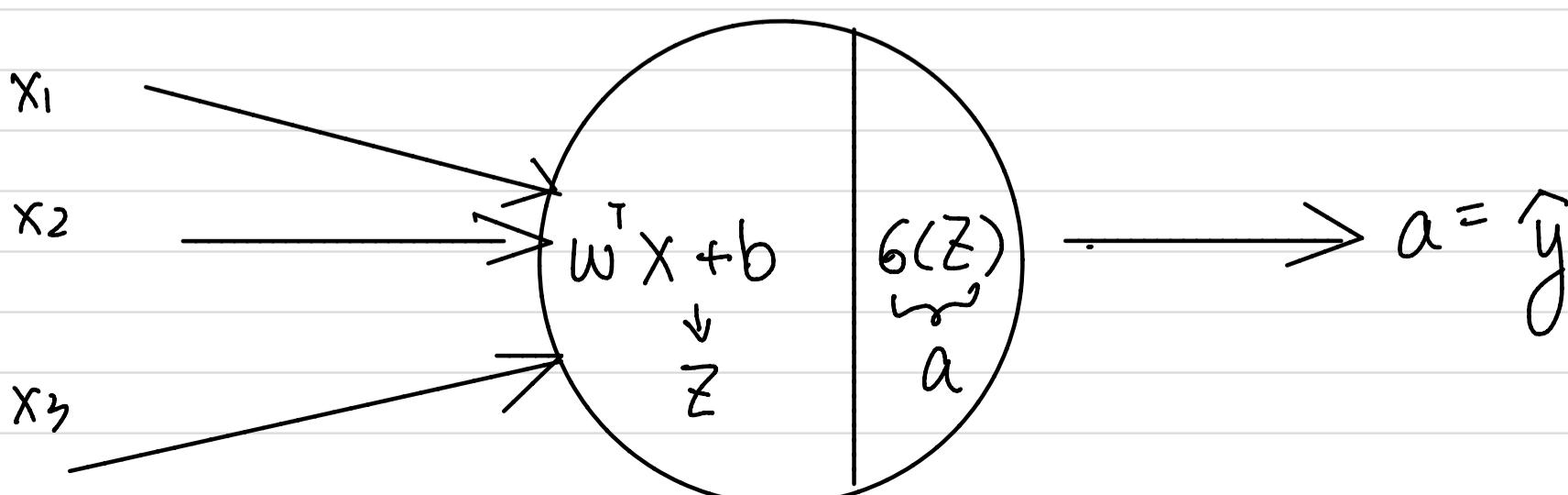
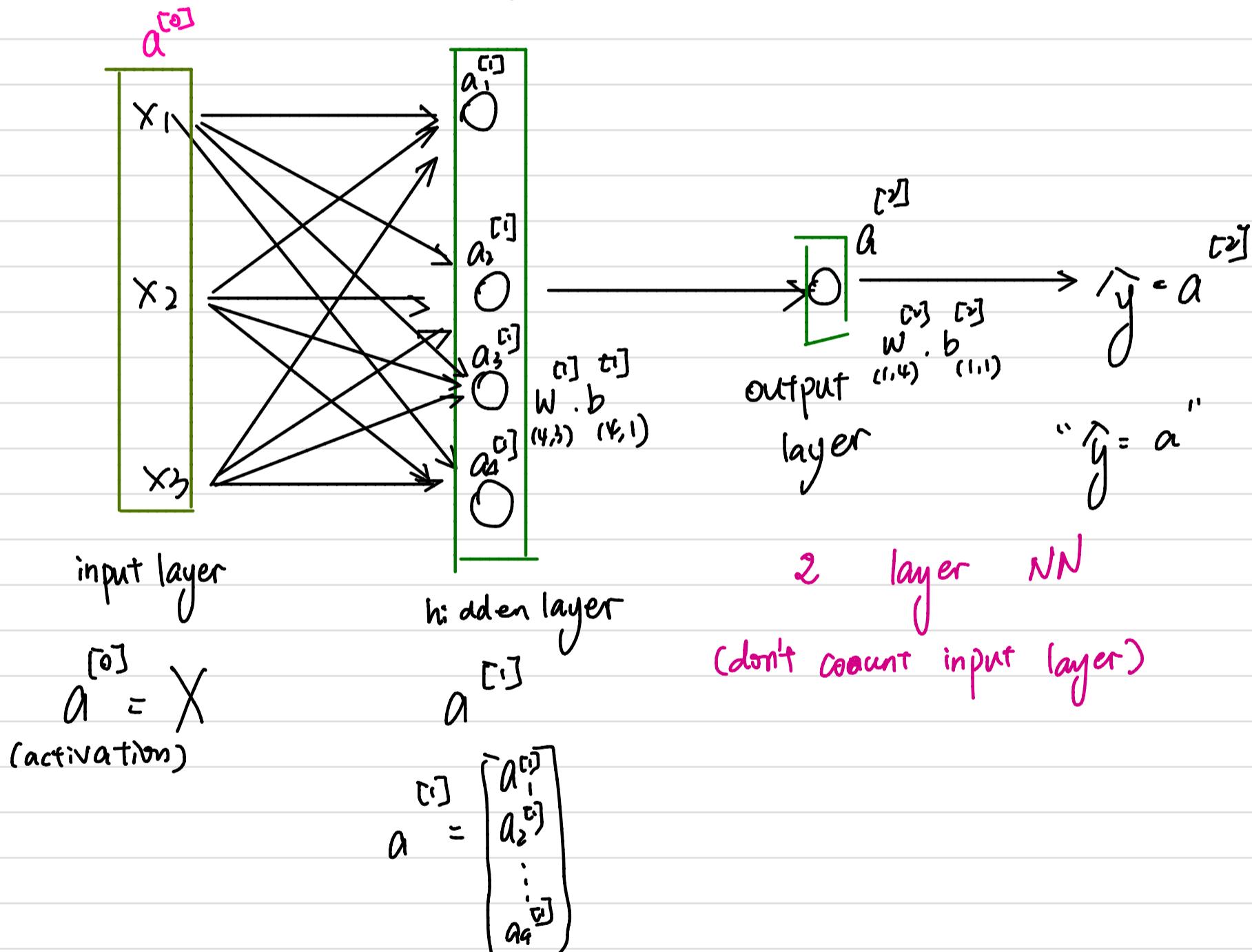
$$z^{[1]} = W^{[1]} X + b^{[1]} \rightarrow a^{[1]} = g(z^{[1]})$$

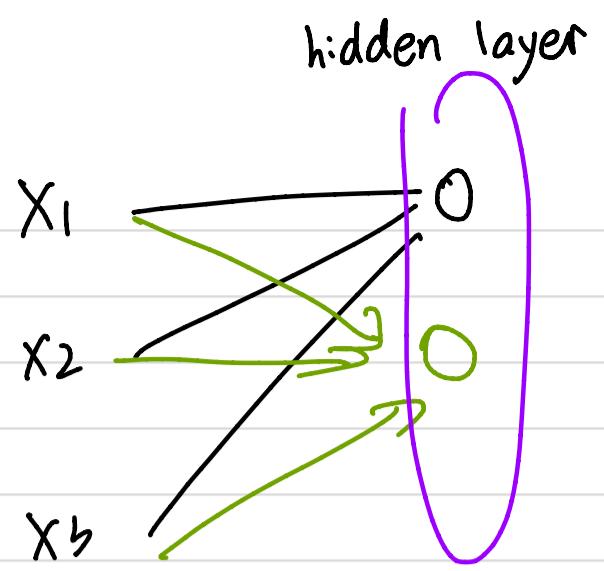
$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = g(z^{[2]})$$

$$\hat{y} = L(a^{[2]}, y)$$

$x^{(1)}$: The first data set (the quantity of one data set is the number of data)

backward calculation is derivative





$$W^T = W^{[2 \times 3]} \\ b = b^{[2 \times 1]}$$

$$z_1^{[1]} = w_1^{[1 \times 3]} x + b_1^{[1 \times 1]} \\ a_1^{[1]} = g(z_1^{[1]})$$

$a_i^{[l]}$ \leftarrow node in layer
 i
 $a_i^{[l]}$ \leftarrow node in layer.

$$z_2^{[1]} = w_2^{[1 \times 3]} x + b_2^{[1 \times 1]} \\ a_2^{[1]} = g(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1 \times 3]} x + b_3^{[1 \times 1]} \\ a_3^{[1]} = g(z_3^{[1]}) \\ z_4^{[1]} = w_4^{[1 \times 3]} x + b_4^{[1 \times 1]} \\ a_4^{[1]} = g(z_4^{[1]})$$

To Matrix:

$$Z^{[4 \times 1]} = \begin{bmatrix} \dots & W_1^{[1 \times 3] T} \\ \dots & W_2^{[1 \times 3] T} \\ \dots & W_3^{[1 \times 3] T} \\ \dots & W_4^{[1 \times 3] T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1 \times 1]} \\ b_2^{[1 \times 1]} \\ b_3^{[1 \times 1]} \\ b_4^{[1 \times 1]} \end{bmatrix} =$$

$$\begin{bmatrix} W_1^{[1 \times 3] T} x + b_1^{[1 \times 1]} \\ W_2^{[1 \times 3] T} x + b_2^{[1 \times 1]} \\ W_3^{[1 \times 3] T} x + b_3^{[1 \times 1]} \\ W_4^{[1 \times 3] T} x + b_4^{[1 \times 1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = g(z^{[1]})$$

$$z^{[2 \times 1]} = w^{[2 \times 3]} a^{[3 \times 1]} + b^{[2 \times 1]}$$

$$a^{[2 \times 1]} = g(z^{[2 \times 1]})$$

$$\begin{aligned}
 x &\rightarrow a^{[2]} = \hat{y} \\
 x^{(1)} &\rightarrow a^{[2]^{(1)}} = \hat{y}^{(1)} \\
 x^{(2)} &\rightarrow a^{[2]^{(2)}} = \hat{y}^{(2)} \\
 &\vdots \\
 x^{(m)} &\rightarrow a^{[2]^{(m)}} = \hat{y}^{(m)}
 \end{aligned}$$

$a^{[2]^{(i)}}$ example i^{th}

layer 2

We define the matrix to be equal to our training examples stacked up in these columns like so. So take the training examples and stack them in columns. So this becomes a $n \times m$ maybe $n \times m$ diminish the matrix

$X = \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 2 & & & & \\ 3 & & & & \\ 4 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix}$

example: 10人測量身高、體重、胖活量.
這個matrix是 3×10 .
↑
attributes
↑
number of people

Use matrix broadcasting.

$$\begin{aligned}
 z^{[1]} &= w^{[1]} X + b^{[1]} \\
 a^{[1]} &= g(z^{[1]}) \\
 z^{[2]} &= w^{[2]} a^{[1]} + b^{[2]} \\
 a^{[2]} &= g(z^{[2]})
 \end{aligned}$$

for loop:

$$\begin{aligned}
 z^{[1]^{(i)}} &= w^{[1]} x^{(i)} + b^{[1]} \\
 a^{[1]^{(i)}} &= g(z^{[1]^{(i)}}) \\
 z^{[2]^{(i)}} &= w^{[2]} a^{[1]^{(i)}} + b^{[2]} \\
 a^{[2]^{(i)}} &= g(z^{[2]^{(i)}})
 \end{aligned}$$

Vectorizing.

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} x_{11} & x_{21} & \dots & x_{m1} \\ x_{12} & & & \\ \vdots & & & \vdots \\ x_{1n_x} & x_{2n_x} & \dots & x_{mn_x} \end{bmatrix} (n_x, m)$$

Capital :

$$\begin{aligned}
 Z^{[1]} &= W^{[1]} X + b^{[1]} \\
 A^{[1]} &= g(Z^{[1]}) \\
 Z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \\
 A^{[2]} &= g(Z^{[2]})
 \end{aligned}$$

$$Z = \begin{bmatrix} \vdots & \vdots & \vdots \\ z^{[1]^{(1)}} & z^{[1]^{(2)}} & \dots & z^{[1]^{(m)}} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

first hidden unit on the first training example
second hidden unit on the first training example

$A^{[1]} = \begin{bmatrix} \vdots & \vdots & \vdots \\ a^{[1]^{(1)}} & a^{[1]^{(2)}} & \dots & a^{[1]^{(m)}} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$

#th training example.

$$z^{1} = w^{[1]} x^{(1)} + b^{[1]}, \quad z^{2} = w^{[2]} x^{(2)} + b^{[2]}, \quad z^{3} = w^{[3]} x^{(3)} + b^{[3]}$$

$$w^{[1]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \vdots \\ \cdot \end{bmatrix}$$

$$w^{[2]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \vdots \\ \cdot \\ b \end{bmatrix}$$

$$w^{[3]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \vdots \\ \cdot \end{bmatrix}$$

$$w^{[1]} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}_{4 \times 3} = \begin{bmatrix} z^{1} & z^{[1](2)} & z^{[1](3)} \\ z^{[2]} & z^{[3]} & z^{[4]} \\ \vdots & \vdots & \vdots \end{bmatrix}_{[1] \times 3}$$

$= z$

Activation functions.

- ① sigmoid
- ② tanh
- ③ ReLU

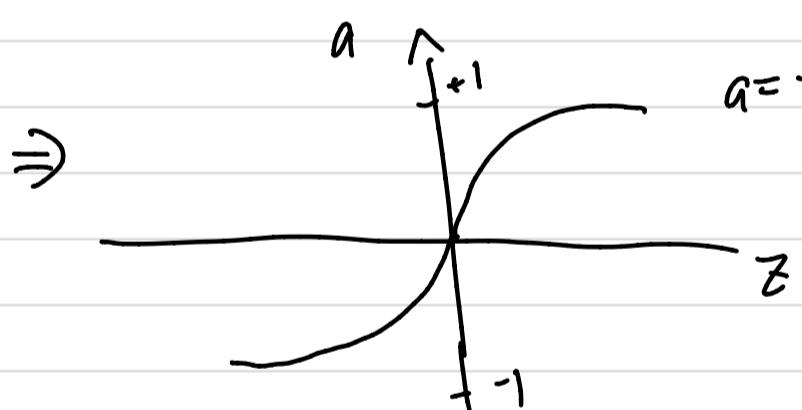
sigmoid



$$a^{[1]} = g(z^{[1]})$$

$$g(z^{[1]})$$

$$a^{[2]} = g(z^{[2]}).$$



$$a = \tanh(z)$$

$$a = \frac{e^{-z} - e^z}{e^{-z} + e^z}$$

(shifted version
of sigmoid).

$$g(z^{[1]}) = \tanh(z^{[1]}).$$

so data has a 0 mean.

$$y \in \{0, 1\}$$

$$1 \quad 1$$

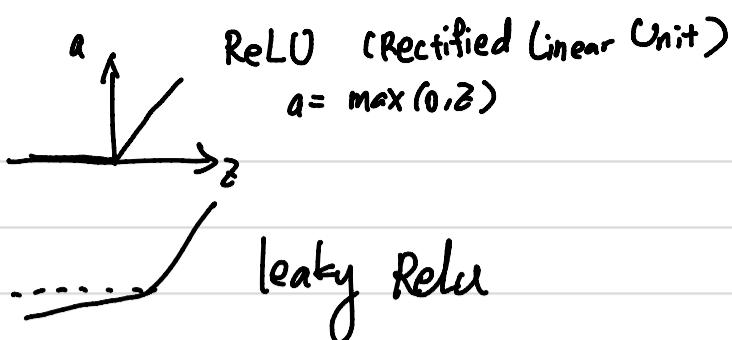
hidden layer

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

$$g^{[1]}(z^{[1]}) = \tanh(z^{[1]}).$$

tanh

output layer: $g^{[2]}(z^{[2]}) = g(z^{[2]})$ so use sigmoid for the output layer.



Pros and cons.

sigmoid : never use it (except for output layer sometimes)

$$\text{we use tanh. } a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

* the most used is ReLU $a = \max(0, z)$

leaky ReLU $a = \max(0.01z, z)$

* Non-linear Regression.

Derivative

① sigmoid

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

$$g'(z) = \frac{d}{dz} g(z) = \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}}\right) = g(z)(1 - g(z))$$

② tanh

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

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

$$= 1 - \tanh^2(z)$$

$$= 1 - a^2$$

③ ReLU

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined} & z = 0 \end{cases}$$

④ Leaky ReLU

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined} & z = 0 \end{cases}$$

Your neural network which only has one hidden layer has parameters:

$$(n^{[1]}, n^{[0]}) \cdot (n^{[1]}, 1) \cdot (n^{[2]}, n^{[1]}) \cdot (n^{[2]}, 1)$$

$n_x = n^{[0]}$ input features. $n^{[1]}$ hidden units $n^{[2]}$ output units = 1

Cost function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$

Gradient Descent

$$\hat{y}^{(i)}, i=1, \dots, m$$

$$\begin{aligned} dW^{[1]} &= \frac{\partial J}{\partial W^{[1]}} \\ W^{[1]} &= W^{[1]} - \alpha dW^{[1]} \quad (\text{update}) \\ b^{[1]} &= b^{[1]} - \alpha db^{[1]} \end{aligned}$$

i don't know

why.

Forward propagation

$$\begin{aligned} Z^{[0]} &= W^{[0]} X + b^{[0]} \\ A^{[0]} &= g^{[0]}(Z^{[0]}) \end{aligned}$$

$$\begin{aligned} Z^{[1]} &= W^{[1]} A^{[0]} + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) = G(Z^{[1]}). \quad \text{here we use sigmoid} \end{aligned}$$

Backward propagation (derivative)

$$\begin{aligned} dz^{[2]} &= A^{[2]} - Y \quad Y = [\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)}] \\ dw^{[2]} &= \frac{1}{m} dz^{[2]} A^{[2]T} \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True}) \end{aligned}$$

$$dz^{[1]} = \underbrace{W^{[1]T} dz^{[2]}}_{(n^{[0]}, m)} * \underbrace{g^{[1]'}(Z^{[1]})}_{(n^{[1]}, 1)}$$

element wise product

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

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Derivative . Backward.

$$\begin{aligned}
 & \text{Forward pass: } z^{[1]} = W^{[1]}X + b^{[1]} \rightarrow a^{[1]} = g(z^{[1]}) \rightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \rightarrow a^{[2]} = g(z^{[2]}) \rightarrow L(a^{[2]}, y) \\
 & \quad \frac{\partial z^{[1]}}{\partial X} = W^{[1]T}, \quad \frac{\partial z^{[2]}}{\partial a^{[1]}} = W^{[2]T} \\
 & \quad * g'(z^{[1]}) \text{ column vector} \rightarrow \frac{\partial w^{[2]}}{\partial b^{[2]}} = \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial b^{[1]}} \rightarrow \frac{\partial w^{[2]}}{\partial b^{[2]}} = \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial b^{[1]}} \text{ row vector}
 \end{aligned}$$

$$\begin{array}{ccc}
 x_1 & 0 & \\
 x_2 & 0 & \\
 x_3 & 0 & \\
 \hline
 n_x = n^{[0]} = 3 & n^{[1]} = 4 & n^{[2]} = 1
 \end{array}$$

(No. of nodes in hidden unit)

$W^{[2]}$ is $n^{[2]} \times n^{[1]}$ dimension

$z^{[2]}, dz^{[2]}$ ($n^{[2]}, 1$) dimension

$z^{[1]}, dz^{[1]}$ ($n^{[1]}, 1$)

dimension

$$\begin{aligned}
 dz^{[2]} &= W^{[2]T} dz^{[1]} * g'(z^{[1]}) \\
 (n^{[2]}, 1) & (n^{[1]}, n^{[2]}) (n^{[2]}, 1) (n^{[2]}, 1).
 \end{aligned}$$

$$dw^{[2]} = dz^{[1]} \cdot W^{[2]T}$$

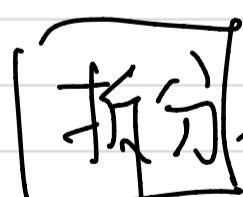
$$db^{[2]} = dz^{[1]}$$

Initialize weights to zero.

Weight: W

建模过程.
每个 cl module 一个.

Random Initialization



拿一个样本. 预测.

-1 layer 把 input tensor \rightarrow output tensor.

样本个数多. 10^7 data 还是 10^7 .

breakdown.

CNN operator. image classification

Question: Does each node (in the same layer) have different W value like:

$$x_1 \quad w^{[r,j]} \quad b^{[r,j]}$$

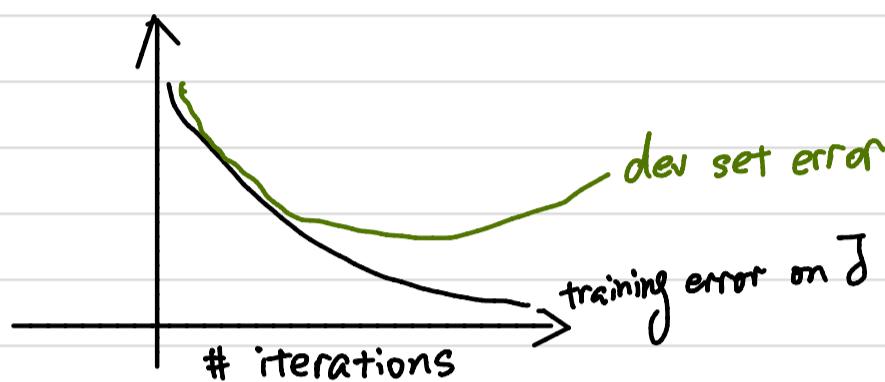
$$O \rightarrow w^{[r,j]}$$

$$x_2 \quad O \rightarrow w^{[r,j]}$$

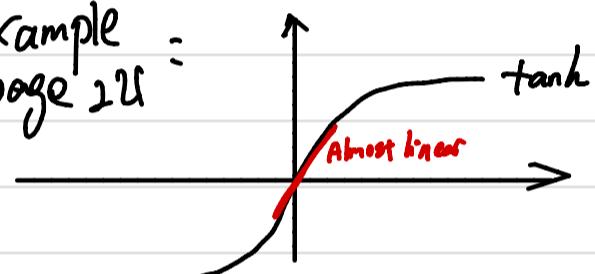
$$O \rightarrow w^{[r,j]}$$

$$x_3 \quad O \rightarrow w^{[r,j]}$$

iteration VS epoch.



Regularization : Example on page 211 =



$$g(z) = \tanh(z)$$

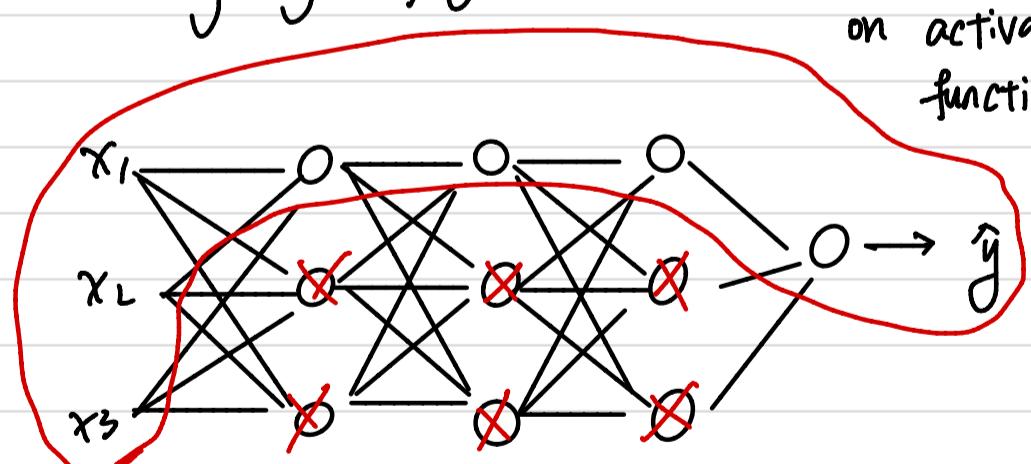
$$\pi \uparrow \quad w^{[l,j]} \downarrow \quad z^{[l,j]} = w^{[l,j]} a^{[l-1]} + b^{[l,j]}$$

\Downarrow
 $g^{[l,j]}(z^{[l,j]})$ will be roughly linear

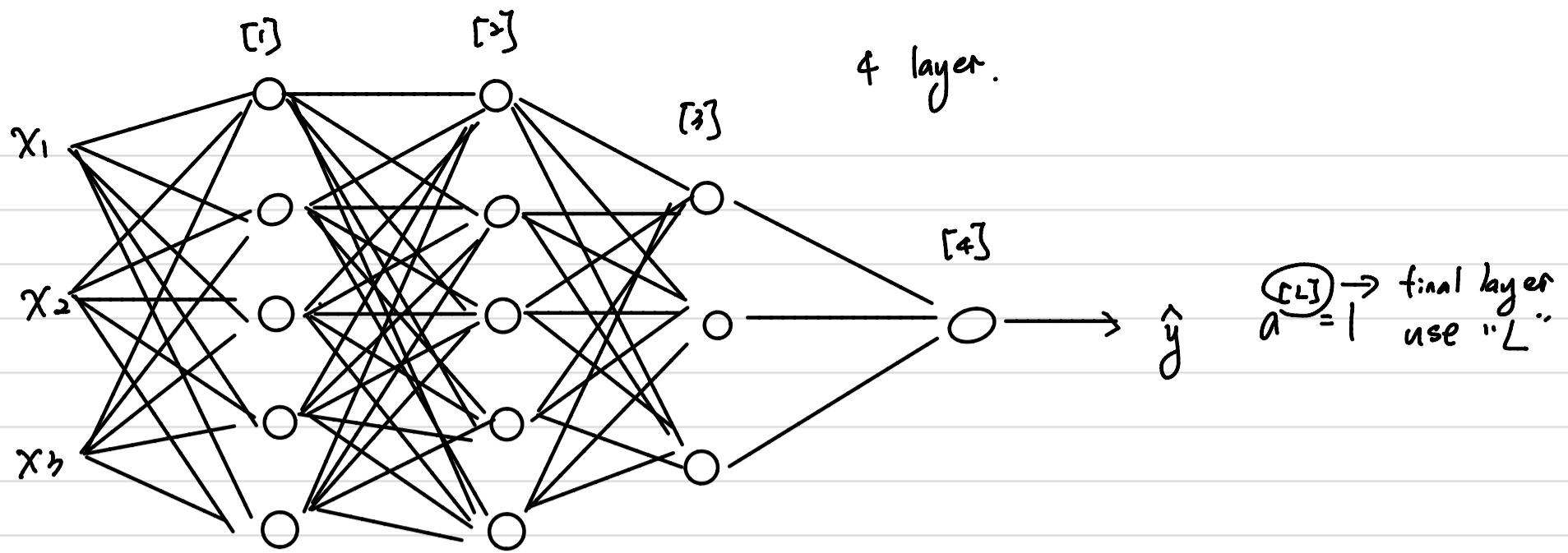
所以說 $w^{[l,j]}$ 越小就越趨近於 linear.

why say roughly linear, doesn't it need to depend on activation function?

Dev data



原理: Zero out units in the layers



$$l = 4$$

$n^{[l]}$ = # units in layer l $n^{[0]}=3, n^{[1]}=5, n^{[2]}=5, n^{[3]}=3, n^{[4]}=1$
 $a^{[l]}$ = activations in layer l ($\text{e.g. sigmoid / tanh / ReLU etc.}$)
 $a^{[l]} = g^{[l]}(z^{[l]})$ ($\text{e.g. sigmoid / tanh / ReLU etc.}$)
 $w^{[l]}$ = weights for $z^{[l]}$
 $b^{[l]}$ = bias for $z^{[l]}$

Forward Propagation
 $X: z^{[0]} = w^{[0]}X + b^{[0]}$
 $a^{[0]} = g^{[0]}(z^{[0]})$

$$z^{[1]} = w^{[1]}a^{[0]} + b^{[1]} \\ a^{[1]} = g^{[1]}(z^{[1]})$$

for loop.

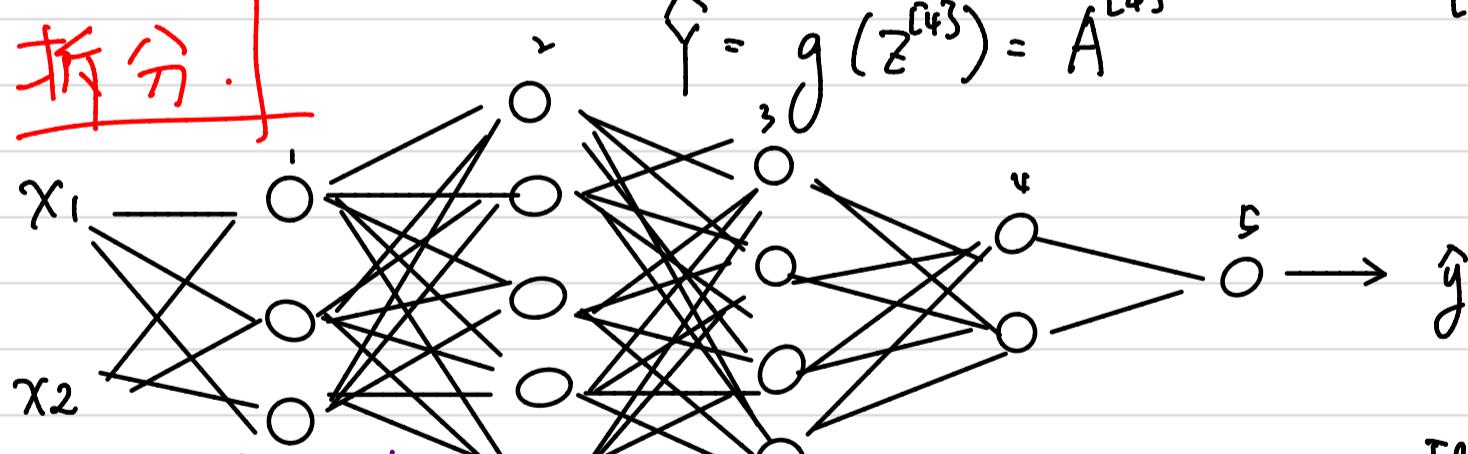
for $l = 1, \dots, 4$ (This for loop is unavoidable).

$$\text{general } z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]}$$

$$z^{[4]} = w^{[4]}a^{[3]} + b^{[4]} \\ a^{[4]} = g^{[4]}(z^{[4]})$$

$$\text{Vector form} \\ z^{[l]} = w^{[l]}A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}(z^{[l]}) \\ z^{[l]} = w^{[l]}A^{[l]} + b^{[l]} \rightarrow z^{[l]} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ z^{[l](1)} & z^{[l](2)} & \dots & z^{[l](m)} \end{bmatrix}$$

拆分



We have two features
Thinking steps color

$$z^{[1]} = w^{[1]} \cdot x_1 + b^{[1]} \\ (3,1) \quad (5,3) \quad (5,1) \\ \vdots \quad \vdots \quad \vdots$$

$$w^{[1]} = (n^{[0]}, n^{[1]}) \\ w^{[2]} = (n^{[1]}, n^{[2]}) \\ z^{[2]} = w^{[2]}a^{[1]} + b^{[2]} \\ (5,1) \quad (5,3) \quad (3,1) \quad (5,1)$$

$$w^{[3]} \quad (3,5) \\ w^{[4]} \quad (2,4)$$

$$\Rightarrow w^{[l]} = (n^{[l]}, n^{[l-1]}) \\ b^{[l]} = (n^{[l]}, 1)$$

$dw^{[l]}$: same as $w^{[l]}$
 $db^{[l]}$: same as $b^{[l]}$

$$z^{[l]} = g^{[l]}(a^{[l]}) : \text{same as } a^{[l]}$$

Then we have more examples. (m examples)

$$(n^{[1]}, 1) \Rightarrow (n^{[1]}, m) \quad \begin{bmatrix} z^{[1]} \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ \vdots & \vdots & & \vdots \end{bmatrix}_{n^{[1]} \times m}$$

features / attributes

All of the training sets stack horizontally. \Rightarrow each row represents the number of X_1, X_2, \dots, X_k (k is the number of features)

$$z^{[1]} = W^{[1]} \cdot X + b^{[1]} \Rightarrow (n^{[1]}, m)$$

$(n^{[1]}, m), (n^{[1]}, n^{[1]}), (n^{[1]}, m), (n^{[1]}, 1)$.

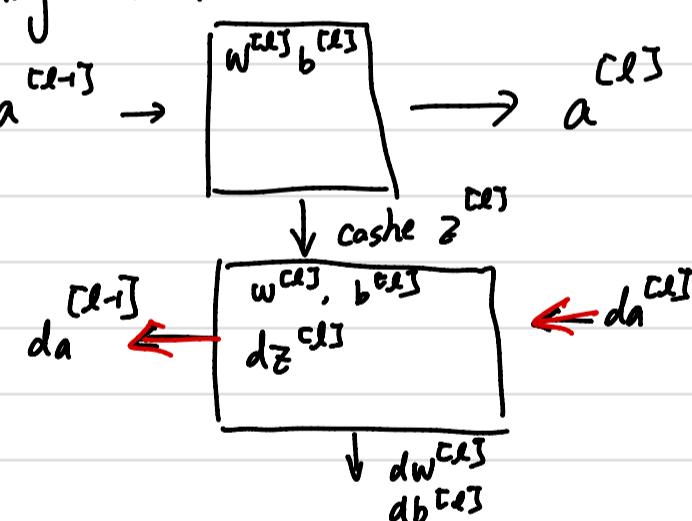
* m is the number of samples
Then we have: $z^{[l]}, A^{[l]}: (n^{[l]}, m)$
 $dz^{[l]}, dA^{[l]}: (n^{[l]}, m)$.

In general, the deep neural network is to 拆分一个图片 to many pixels, for example first layer is to detect edges, then second hidden layer is to detect the facial features like eyes, noses.... then put the third layer us to put them together.,.

Backward Propagation :

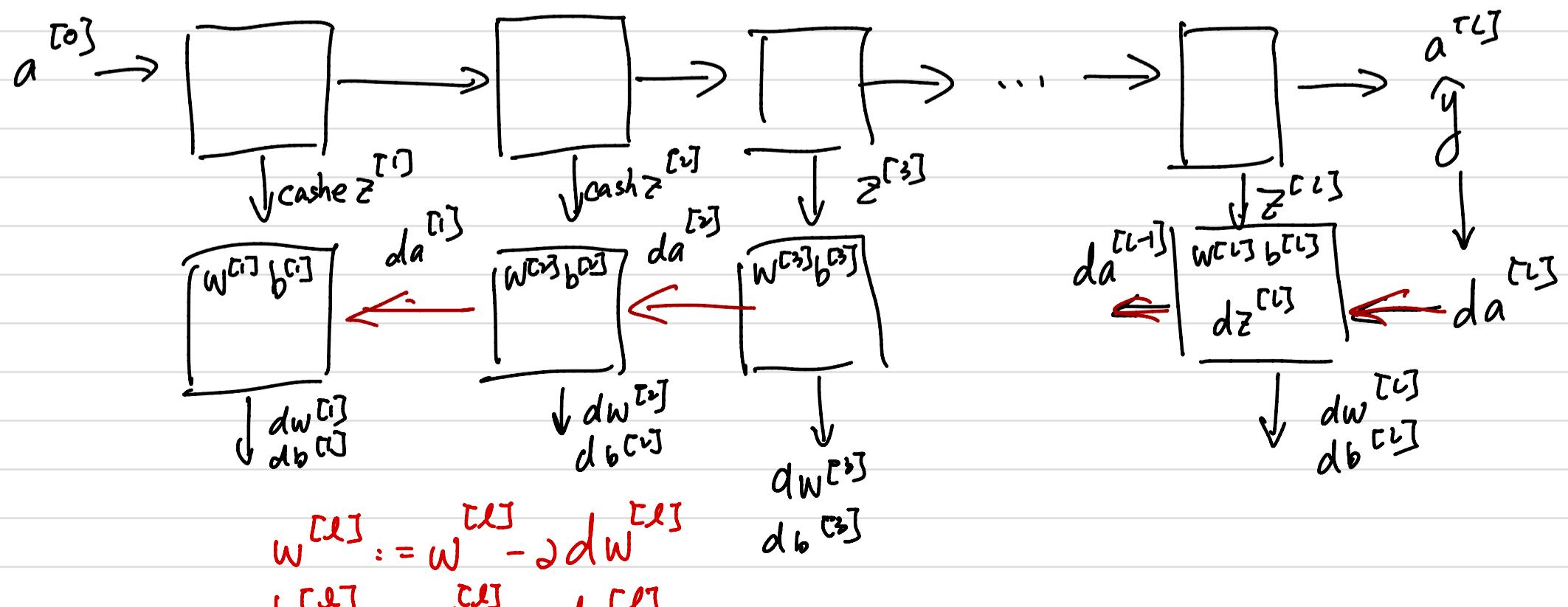
Input $da^{[l]}$ output $da^{[l-1]}$
cashe $z^{[l]}$ $dw^{[l]}$
 $db^{[l]}$

layer l .



What does cashe mean here: cashe is a widely used method for storing information so that it can be later accessed much more quickly.

Forward and Backward.



$$a = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, a \cdot b = a_1 b_1 + a_2 b_2$$

The difference of cross product and dot product

① cross → 结果还是 matrix

② dot → 结果为数字. $U \cdot V = U_1 V_1 + U_2 V_2 + U_3 V_3$

Input $da^{[L]}$

Output $da^{[L-1]}, dw^{[L]}, db^{[L]}$

$$\frac{\partial z^{[L]}}{\partial a^{[L-1]}} = \frac{\partial a^{[L]}}{\partial z^{[L]}} \cdot g^{[L]'}(z^{[L]}) \quad \text{cross product}$$

$$\frac{\partial w^{[L]}}{\partial z^{[L]}} = \frac{\partial z^{[L]}}{\partial z^{[L]}} \cdot a^{[L-1]} \quad \text{dot product}$$

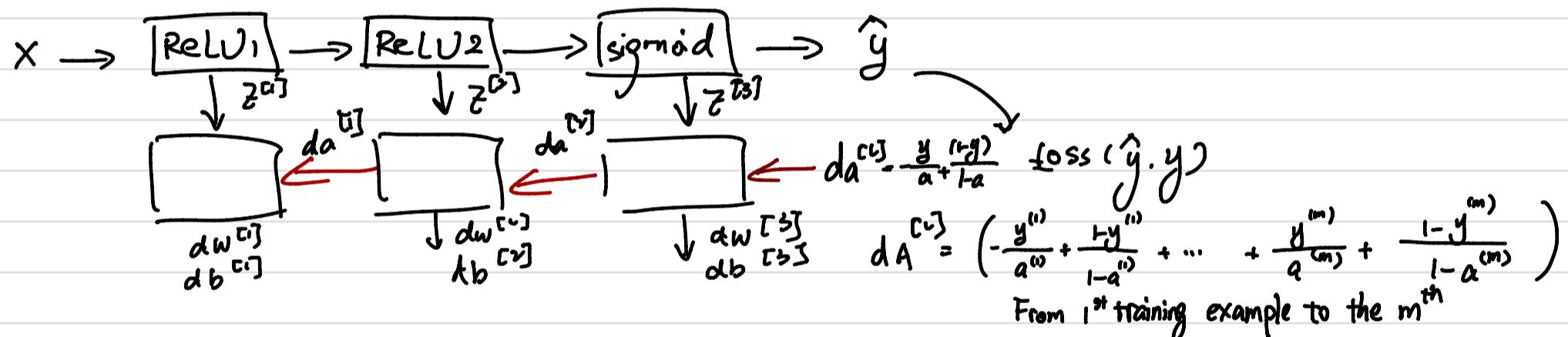
$$\frac{\partial a^{[L-1]}}{\partial z^{[L]}} = W^{[L+1]T} \cdot dz^{[L+1]} + g^{[L]'}(z^{[L]})$$

$$\text{Vectorize : } dz^{[L]} = dA^{[L]} \times g^{[L]'}(z^{[L]})$$

$$dW^{[L]} = \frac{1}{m} dz^{[L]} \cdot A^{[L-1]T}$$

$$db^{[L]} = \frac{1}{m} \text{np.sum}(dz^{[L]}), \text{ axis=1, keepdims=True}$$

$$dA^{[L-1]} = W^{[L]T} \cdot dz^{[L]}$$



Hyperparameter

parameters : $W^{[i]}, b^{[i]}, \dots$

Hyperparameter : learning rate α

iterations

hidden layers L

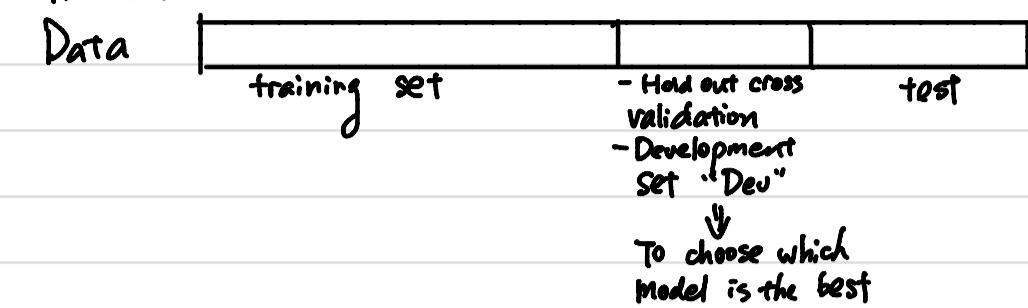
hidden units $n^{[1]}, n^{[2]}, \dots$

choice of activation function

parameters control $W^{[i]}, b^{[i]}$.

The Second Course

Train / dev / test sets.



Prev era. 70/30 60/20/20.

Modern data era. 1,000,000 Dev become smaller

98% / 1% / 1%

99.5% / 0.4% / 0.1%

Mismatched train / test distribution

Training set: cat pics from web

Dev/test sets: cat pics from App. 不清楚.

* Not having a test set might be okay. (only dev set.)

\Rightarrow Make sure that dev and test sets from same distribution

Bias and Variance

	cat	not cat	
cat classification	$y=1$	$y=0$	
Train set error	1%	65%	15%
Dev set error	11%	16%	30%
high variance overfitting	high bias (cat doing well on training set)	high bias high variance	low bias low variance.
	underfitting		

Human 0%.

Optimal error (Bayes) : 0%

Blurry pics. \Rightarrow Bayes error \uparrow

Bias Problem: How well your training data is

Variance Pro: training \rightarrow test. how bad

High bias: The model does not fit well

High var: The model overfits the data.

Thinking:

- High bias? \rightarrow \uparrow
- Bigger network. (* Always reduce bias without hurting variance)
 - NN larger

High variance?

(dev. set performance). \rightarrow

$\downarrow N$

Done

\uparrow

bias)

- More Data (* Always reduce var without hurting bias)
- Regularization

"Bias-variance tradeoff"

λ : weight $\frac{1}{2} \cdot \lambda$, prevent overfitting problem

Logistic Regression

$$\min J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 \quad (w \in \mathbb{R}^n, b \in \mathbb{R})$$

L2 regularization: $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$

L1 regularization: $\frac{\lambda}{2m} \sum_{i=1}^m |w_i| = \frac{\lambda}{2m} \|w\|_1$, (w will be sparse).

In python: λ -Lambda.
So we use lambda for this function.

Neural Network

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2 \quad - L2 \text{ regularization}$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n (w_{ij}^{[l]})^2, \quad W: (n^{[l]}, n^{[l-1]})$$

softmax

"frobenius norm"

Gradient,

$$dw^{[l]} = (\text{from before}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

$$\Rightarrow w^{[l]} := w^{[l]} - \alpha (\text{from before} + \frac{\lambda}{m} w^{[l]})$$

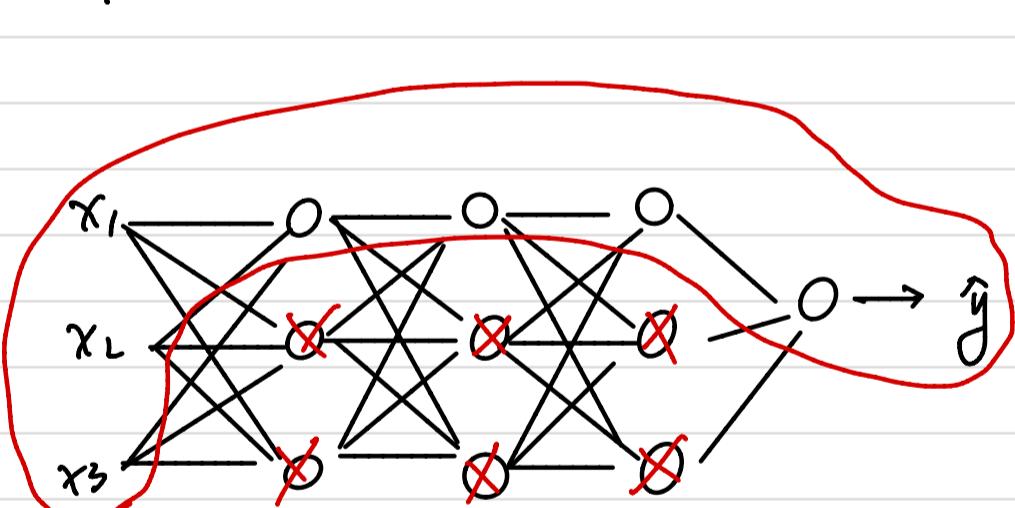
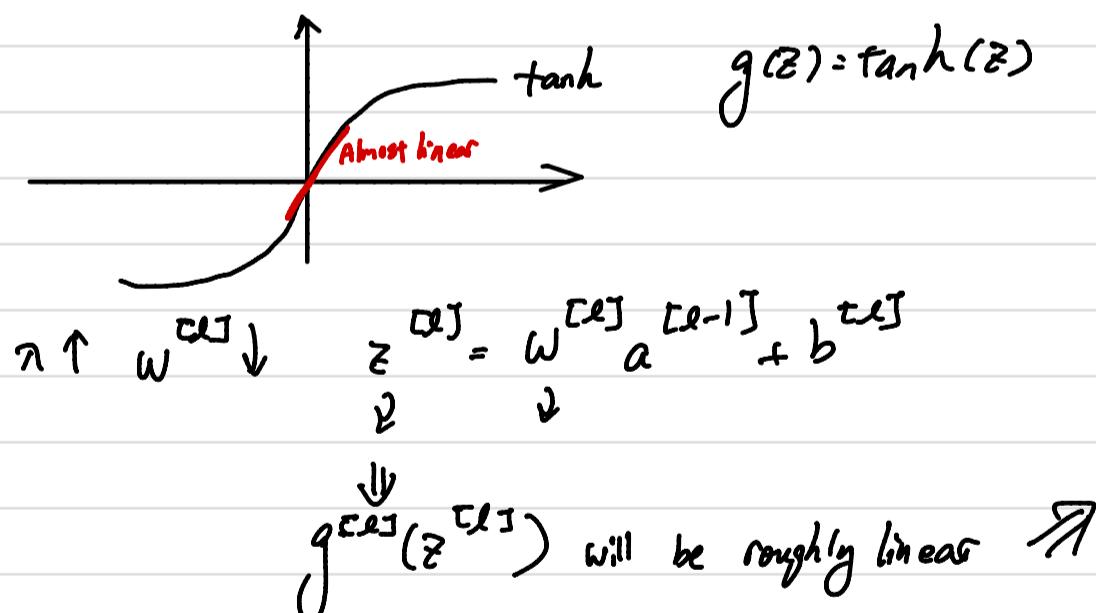
$$= w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from before})$$

"weight decay" $= (1 - \frac{\alpha \lambda}{m}) w^{[l]} - \alpha (\text{from before})$

Why regularization prevent overfitting?

$$J(w^{[0]}, b^{[0]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

If λ is really big, and we want $J(w^{[0]}, b^{[0]})$ to be small
then $w^{[l]} (\approx 0)$, which will leads to a simpler NN (neural network) \Rightarrow overfitting



Dropout (drop some units in different layers) "Inverted Dropout"

layer: $l=3$ keep-prob = 0.8 (the prob that nodes will be kept)

$d_3 = \text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keep-prob}$

random matrix

activation of layer 3.

$a_3 = \text{np.multiply}(a_3, d_3)$ $a_3 * = d_3$

$a_3 / = \text{keep-prob}$ e.g. If we have 50 units in 3rd layer $\because \text{keep-prob} = 0.8 \Rightarrow 20\% \text{ zero}$

then 10 units will be zero out

$$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$$

reduced by 20%

Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights.

Each layer can have different keep-prob. you can also set keep-prob=1. means no units need to be get rid of.



Other organization.

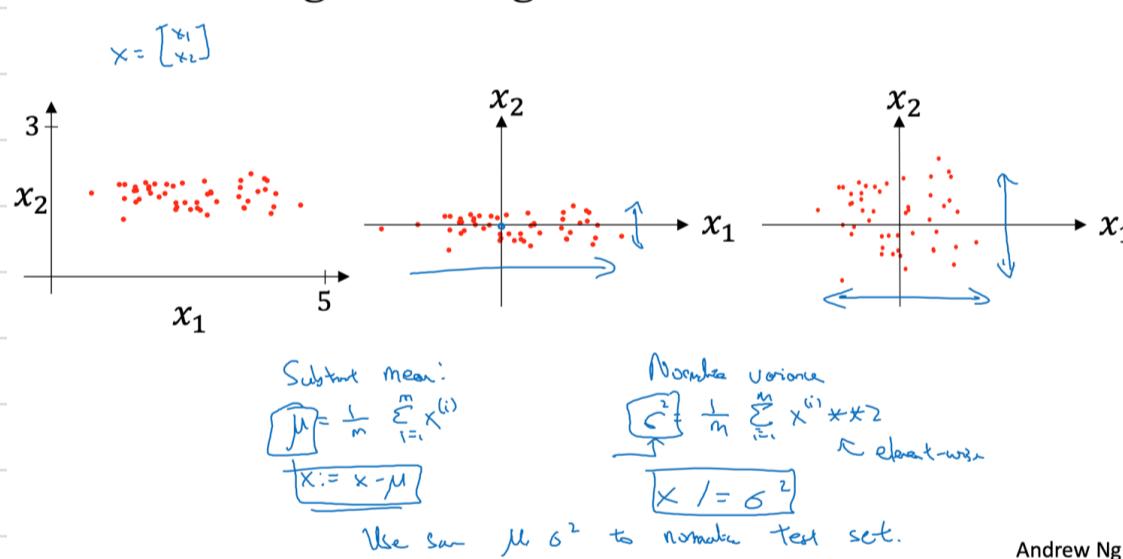
Optimize cost function J : $\Rightarrow J(w, b)$

- Gradient,

Not overfitting:

- Regularization $\|w\|_F^2$.

Normalizing training sets



Vanishing/exploding gradients

