

Deep learning (course)

- ## Deep learning (cont'd)

Input (x)	Output (y)	Application
Home features	Price	Real estate (Standard NN)
Image	object (1, ..., 1000)	Photo tagging (CNN)
Audio	Text transcript	Speech recognition (Recurrent NN)
Image, Radar info	Position of other cars	Autonomous driving custom / hybrid

Input feature vector

- ## * Some notations

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

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

$$n = n_{\text{train}}, \quad m = m_{\text{test}}$$

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

$\leftarrow m$

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

X .shape = (n_x, m)

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

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

y .shape = $(1, m)$

* Logistic Regression Cost function

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

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, count $\hat{y}^{(i)} \approx y^{(i)}$.

$\hat{y}^{(i)}$ ↓
 over output true label

Loss (error) function :-

$$l(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

FOR single training example

If $y=1$: $l(\hat{y}, y) = -\log \hat{y} \leftarrow$ want $\log \hat{y}$ large,
want \hat{y} large
{ for small loss }

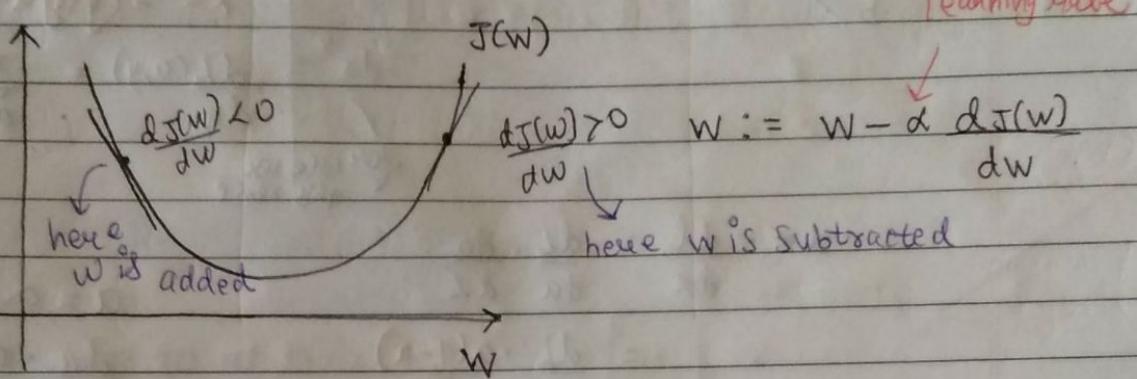
If $y=0$: $l(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $(1-\hat{y})$ to be
large so,
 \hat{y} to be small
{ for small loss }

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

For entire
training set

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

* Gradient descent



* Computation graph

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

$\underbrace{a+bc}_v$

$$U = bc$$

$$V = a+U$$

$$J = 3V$$

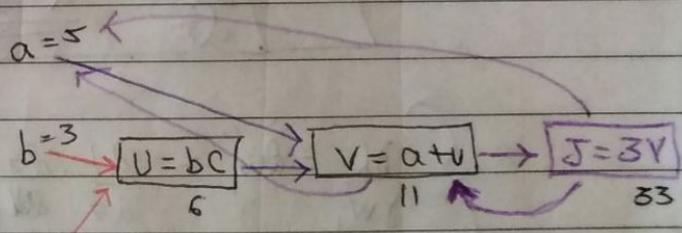
$$\frac{dJ}{dV} = 3$$

$$\frac{dV}{da}$$

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

$$= 3 \times 1$$

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



$$J = 3V$$

$$V = 11 \rightarrow 11.001$$

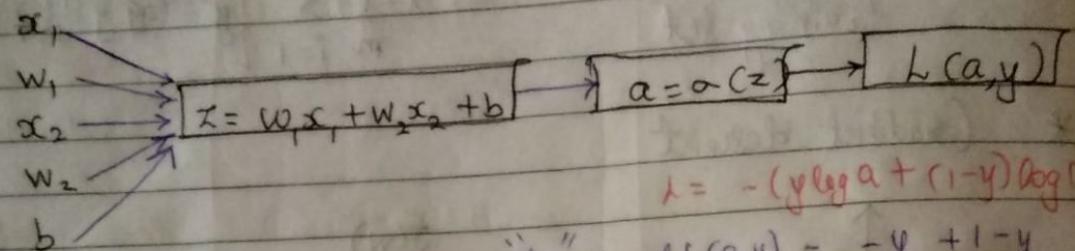
$$J = 33 \rightarrow 33.003$$

$$\{ a = 5 \rightarrow 5.001$$

$$\{ V = 11 \rightarrow 11.001$$

$$\{ J = 33 \rightarrow 33.003$$

* Logistic regression derivatives.



$$\frac{da}{dz} = \frac{dL(a, y)}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$$

for code use
this name

$$\frac{dz}{dx} = \frac{dl}{dz} = \frac{dl}{da} \cdot \frac{da}{dz}$$

$$a = a(z)$$

$$= \frac{dl}{da} \cdot a(1-a)$$

$$\frac{da}{dz}$$

$$\frac{dz}{dx} = \frac{dl}{dz} = \begin{bmatrix} -y & 1-y \\ a & 1-a \end{bmatrix} a(1-a)$$

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

$$= \begin{bmatrix} -y + ay + a - y \\ a - a^2 \end{bmatrix}$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

$$\frac{dz}{dx} = \boxed{\frac{dl}{dz} = \frac{a-y}{a}}$$

$$a(z) = 1$$

$$1+e^{-z}$$

$$\frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial z} = \frac{dl}{dz} \cdot \frac{dz}{\partial w_1} = \boxed{dz \times \text{some } x_1}$$

$$= a^2(z) \cdot e^{-z}$$

$$= a^2(z) \cdot (1-a(z))$$

$$a(z)$$

$$dw_1 = x_1 dz$$

$$db = dz$$

$$= a(z)(1-a(z))$$

$$= a(1-a)$$

$$w_1 := w_1 - \lambda dw_1$$

$$w_2 := w_2 - \lambda dw_2$$

$$b := b - \lambda db$$

* Logistic regression on m examples.

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

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b) \quad \frac{dw_1}{db}, \frac{dw_2}{db}$$

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

$$J = 0; dw_1 = 0; dw_2 = 0; db = 0$$

For $i=1$ to m

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

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

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

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

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

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

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

$\uparrow n=2$ (features)

$$J / i = m$$

$$dw_1 / i = m$$

$$db / i = m$$

$$dw_2 / i = m$$

$$dw_1 = \frac{\partial J}{\partial w_1}, \text{ if for } dw_2 \text{ & } db$$

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

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

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

* Vectorization [To avoid ~~for loops~~]
 (less efficient for bigger data sets)

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

Non-vectorized

$$z = 0$$

For i in range (n_x):
 $z += w^T[i] * x[i]$

$$z += b$$

vectorized (Takes very less time)

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

* Example :-

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

If you want to calculate exponential of each element of v

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

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

For i in range (n):

$$U[i] = \text{math.exp}[v[i]]$$

[non-vectorized] ↴

Now, vectorized

import numpy as np

$$U = np.exp(v)$$

* we use vectorization in our code to get rid of 1 for loop.

$$J = 0, \boxed{dw_1 = 0}, \boxed{dw_2 = 0}, db = 0$$

for $i = 1 \text{ to } m:$

$$dw = np.zeros((n_x, 1))$$

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

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

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

$$dz^{(i)} = a^{(i)} [1 - a^{(i)}]$$

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \right\}$$

$$dw += x^{(i)} dz^{(i)}$$

$$J = J/m, \boxed{dw_1 = dw_1/m}, \boxed{dw_2 = dw_2/m}, \boxed{db = db/m}$$

$$dw = m$$

* Vectorizing logistic regression

$$\begin{aligned} z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b & z^{(3)} &= w^T x^{(3)} + b \\ a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

$$X = \begin{bmatrix} & & \\ \vdots & \vdots & \\ x^{(1)} & x^{(2)} & \dots x^{(m)} \\ & \vdots & \vdots \end{bmatrix} \quad (n_x, m)$$

$$[z^{(1)} \ z^{(2)} \ z^{(3)} \dots z^{(m)}] = w^T \cdot X + [b \ b \dots b]_{1 \times m}$$

$$= \begin{bmatrix} w^{(1)} \text{ (row vector)} \\ w^{(2)} \\ \vdots \\ w^{(m)} \end{bmatrix} \begin{bmatrix} x^{(1)} \ x^{(2)} \dots x^{(m)} \end{bmatrix}$$

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

$$Z = [z^{(1)} \ z^{(2)} \ z^{(3)} \dots z^{(m)}] = [w^T x^{(1)} + b \ w^T x^{(2)} + b \ \dots \ w^T x^{(m)} + b]$$

$$Z = w^T \cdot \text{dot}(w, X) + b$$

[Broadcasting in Python]

$$A = [a^{(1)} \ a^{(2)} \ a^{(3)} \dots \ a^{(m)}] = a(\vec{z})$$

* Vectorizing gradient computation

$$d\vec{z} = [d\vec{z}^{(1)} \ d\vec{z}^{(2)} \ \dots \ d\vec{z}^{(m)}]$$

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

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

$$d\vec{z} = A - y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$

$$dw = 0$$

$$dw + = x^{(1)} d\vec{z}^{(1)}$$

$$dw + = x^{(2)} d\vec{z}^{(2)}$$

\vdots

$$dw / = m$$

we got
end of this
for loop

$$db = 0$$

$$db + = d\vec{z}^{(1)}$$

$$db + = d\vec{z}^{(2)}$$

\vdots

$$db / = m$$

$$db / = m$$

For this,

$$db = \frac{1}{m} \sum_{i=1}^m d\vec{z}^{(i)}$$

$$\boxed{db = \frac{1}{m} np.sum(d\vec{z})}$$

$$\boxed{dw = \frac{1}{m} X d\vec{z}^T}$$

$$= \frac{1}{m} \left[\begin{array}{c|c|c|c} 1 & \dots & 1 \\ x^{(1)} & \dots & x^{(m)} \end{array} \right] \left[\begin{array}{c} d\vec{z}^{(1)} \\ \vdots \\ d\vec{z}^{(m)} \end{array} \right]$$

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

Modified logistic regression

For iteration in range (1000):

$$\hat{z} = w^T x + b$$

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

$$A = \alpha(\hat{z})$$

$$d\hat{z} = A - Y$$

$$dW = \frac{1}{m} X d\hat{z}^T$$

$$db = \frac{1}{m} np \cdot \text{sum}(d\hat{z})$$

$$w = w - \lambda dW$$

$$b = b - \lambda db$$

Broadcasting example

Calories from Carbs, proteins, Fats in 100g of different foods:

	Apples	Beef	Eggs	Potatoes	
Carbs	56.0	0.0	4.4	68.0	
Protein	1.2	104.0	52.0	8.0	
Fat	1.8	135.0	99.0	0.9	

$= A \quad (3, 4)$

Suppose you want to calculate % of calories from carb, protein, fat.

$$① \text{cal} = A \cdot \text{Sum}(\text{axis}=0)$$

means it will sum
i.e. a row vector

$$② \text{percentage} = 100 * \text{a} / (\text{cal})$$

If $(\text{axis}=1)$ it will sum

i.e. Column vector

$\{ \}$

* Broadcasting example

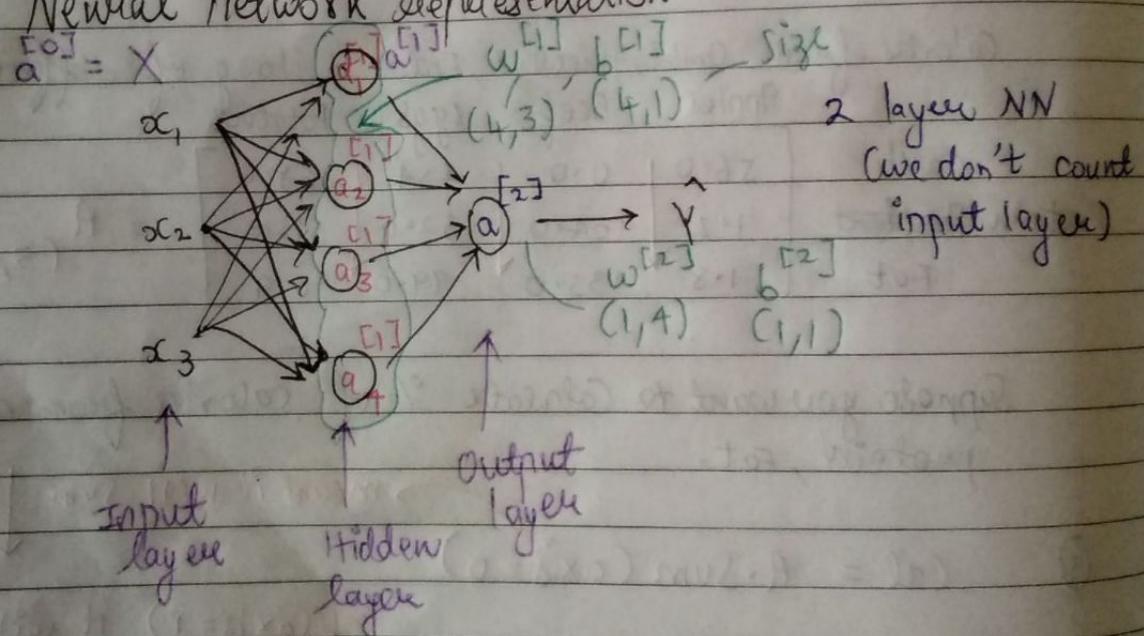
$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

\downarrow

$$\begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

* Neural network representation



$$\Sigma_1 = w_1^{[1]T} x + b_1^{[1]}$$

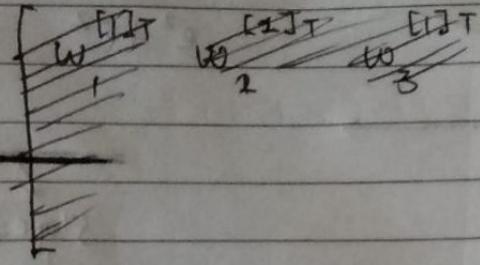
$$a_1^{[1]} = \sigma[\Sigma_1]$$

$$\Sigma_2 = w^{[2]T} a^{[1]} + b^{[2]}$$

$$a_2^{[2]} = \sigma[\Sigma_2]$$

$a_i^{[l]}$ — layer

i — node in layer



$$\begin{bmatrix} \xrightarrow{\quad w_1^{[1]T} \quad} \\ \xrightarrow{\quad w_2^{[1]T} \quad} \\ \xrightarrow{\quad w_3^{[1]T} \quad} \\ \xrightarrow{\quad w_4^{[1]T} \quad} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} =$$

(weight betw 0th & 1st
(layer))

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}_{4 \times 1}$$

* Vectorizing across multiple examples.

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix}_{(n_x, m)} \quad z^{[1]} = \begin{bmatrix} z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}$$

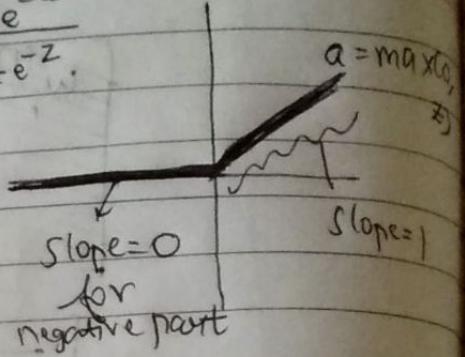
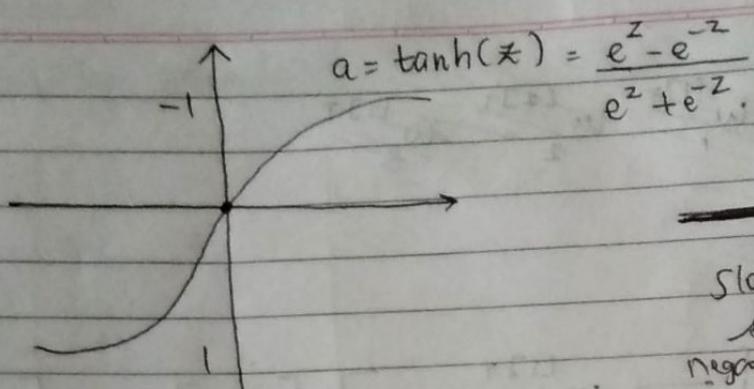
layer no training example no

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

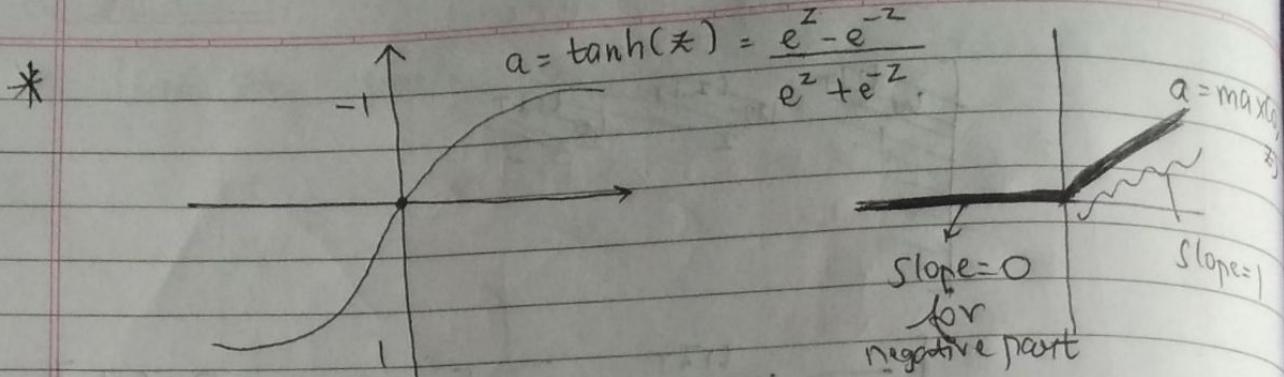
hidden unit

Training examples

*



- * For output layer you can use Sigmoid fn
- * For hidden layer you can use Relu.



- * For output layer you can use Sigmoid fⁿ.
- * For hidden layer you can use ReLU.

* Gradient descent for neural networks.

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ $\rightarrow (n^{[2]}, 1)$
 $(n^{[1]}, n^{[0]}) \quad \quad \quad (n^{[1]}, 1) \quad \quad (n^{[2]}, n^{[1]})$

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

Formulas for Computing derivatives

Forward propagation:

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

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

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

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

Back propagation :-

$$d\hat{z}^{[2]} = A^{[2]} - Y$$

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

$$dW^{[2]} = \frac{1}{m} d\hat{z}^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(d\hat{z}^{[2]}, axis=1, keepdims=True)$$

prevent python
from showing
(n^[2],)

$$d\hat{z}^{[1]} = W^{[2]T} \underbrace{d\hat{z}^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]'}(\hat{z}^{[1]})}_{(n^{[1]}, m)}$$

element wise product.

$$dW^{[1]} = \frac{1}{m} d\hat{z}^{[1]} X^T$$

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

would prevent python
from showing
(n^[1],)

For proof of :-

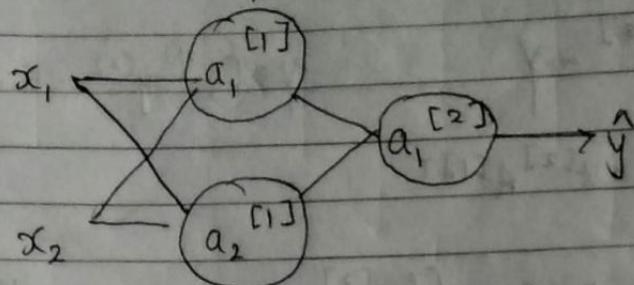
$$d\hat{z}^{[1]} = W^{[2]T} d\hat{z}^{[2]} * g^{[1]'}(\hat{z}^{[1]})$$

$$X^{[1]} \rightarrow a^{[1]} \rightarrow Z^{[2]} \rightarrow a^{[2]} \rightarrow L$$

Now,

$$d\hat{z}^{[1]} = \frac{dL}{d\hat{z}^{[1]}} = \frac{dL}{d\hat{z}^{[2]}} \cdot \frac{d\hat{z}^{[2]}}{da^{[2]}} \cdot \frac{da^{[2]}}{dZ^{[2]}}$$

* Random initialization



$$w^{[1]} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{[1]} = \text{np.zeros}(2, 1)$$

similarly,

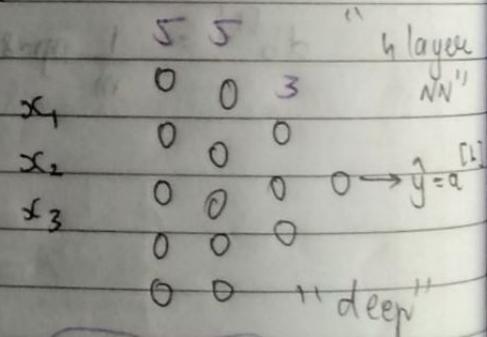
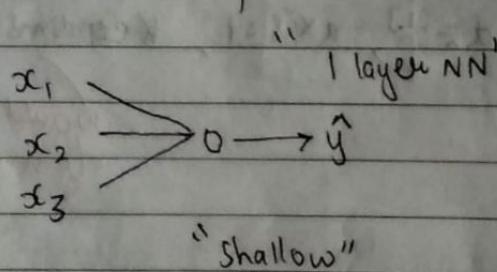
$$w^{[2]} = \dots$$

$$b^{[2]} = \dots$$

To make
random no's
smaller

(as we are
using Sigmoid
& other activ
fn)

* What is a deep neural network?



$$L = 4$$

$n^{(l)}$ = no of units or
nodes in layer

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]}$$

$$n^{[0]} = n = 3$$

$$a^{[l]} = g^{[l]}(\mathbf{x}^{[l]})$$

$$w^{[l]} = \text{weights for } x^{[l]}$$

* Forward propagation in a deep network.

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

$$w^{[l]} = (n^{[l]}, n^{[l-1]}) = dw^{[l]}$$

↓ dimension of $w^{[l]}$ ↑

$$b^{[l]} = (n^{[l]}, 1) = db^{[l]}$$

* Vectorized implementation

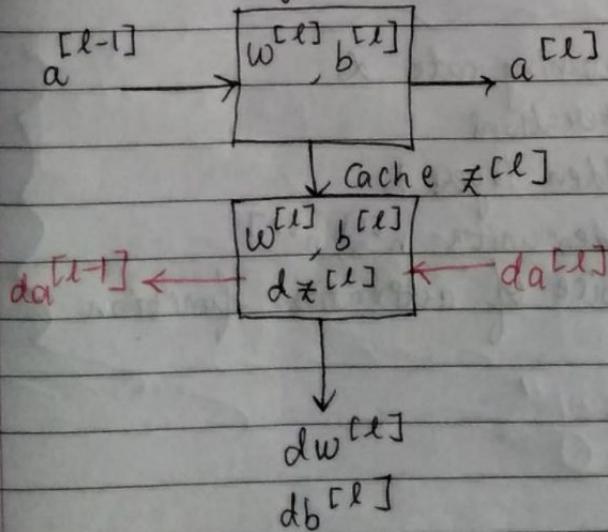
$$z^{[l]}, a^{[l]} = (n^{[l]}, 1)$$

$$z^{[l]}, A^{[l]} = (n^{[l]}, m) = dz^{[l]} = dA^{[l]}$$

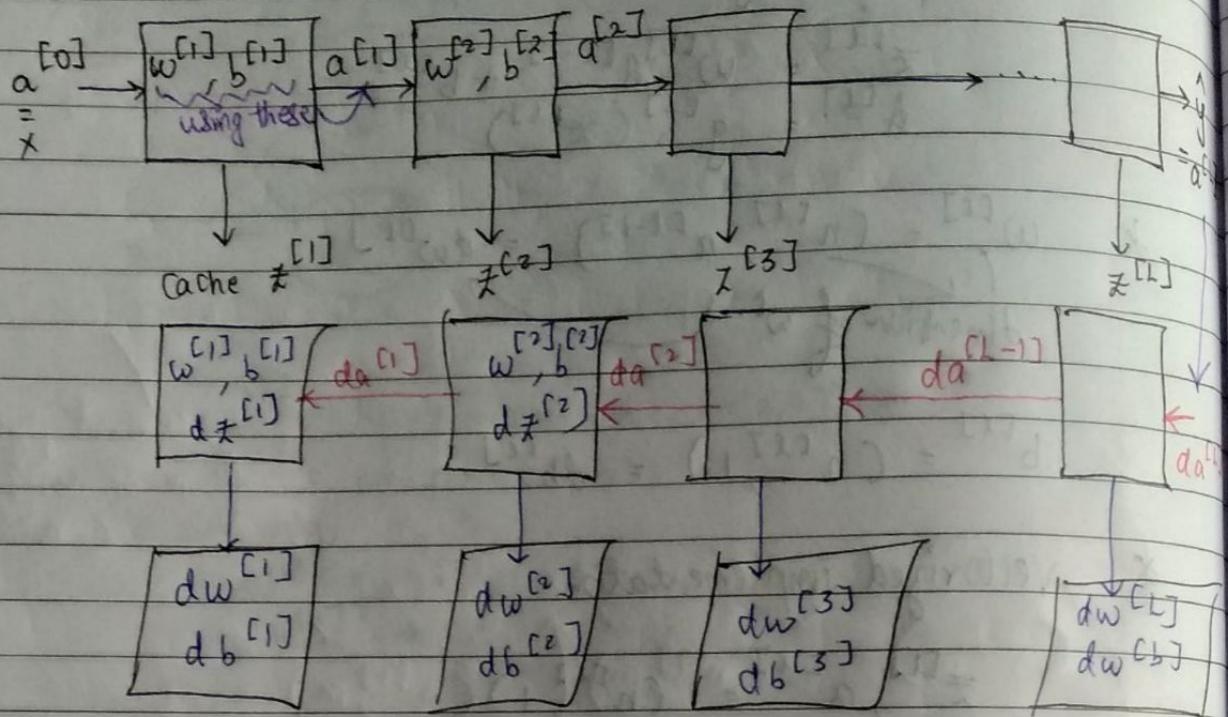
when $l=0$,

$$A^{[0]} = x = (n^{[0]}, m)$$

* Forward & backward functions.
layer l



One iteration of gradient descent



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

$$b^{[l]} := b^{[l]} - \lambda db^{[l]}$$

* What are hyperparameters :-

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$

hyperparameters:

- learning rate λ
- Iterations
- hidden layers L
- hidden units $n^{[1]}, n^{[2]}, \dots$
- choice of activation function

Data

training set

for training
algorithms.

Hold-out
cross validation

development set

"dev"

to find which
model performs
best

test

to evaluate
the model
that we
selected
in dev.

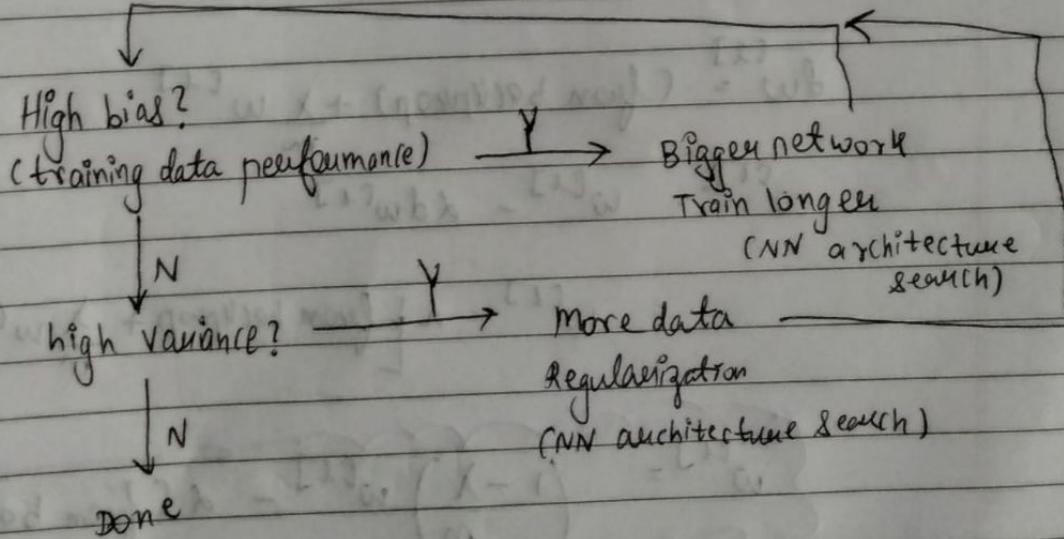
→ Make sure dev & test come from same distribution.

Bias and Variance

Train set error: 1%	15%	15%	0.5%
---------------------	-----	-----	------

Dev set error: 11%	16%	30%	1%
--------------------	-----	-----	----

high bias high variance both high low bias
 Variance bias high low variance.



Logistic regression

Regularization parameter

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$\|w\|_1 = \frac{\lambda}{2m} \sum_{i=1}^{n_x} |w_i| = \frac{\lambda}{2m} \|w\|_1 \quad (w \text{ will be sparse})$$

it will have
lot of zeros
in it)

For Neural network :-

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

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

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

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

$$= w^{[l]} - \lambda \left[\text{from backprop} + \frac{\lambda}{m} w^{[l]} \right]$$

$$w^{[l]} = \left(1 - \frac{\lambda}{m}\right) w^{[l]} - \alpha [\text{from backprop}]$$

"Weight decay"

2)

Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$

$d_3 = \text{np.random.rand}(a_3.shape[0], a_3.shape[1]) <$
 dropout vector for say Keep-prob
 layer 3

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

$\underline{a_3} = 0.8 \text{ (Keep-prob)}$

Let's say you have

50 units \rightarrow 10 units are shut off

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

To maintain original value of $x^{[4]}$

reduced by 20%

we divide a^3 .

* Try to keep-prob high where you think there is not much need of removing neurons. ex: For $(4, 3)$ matrix (weight)

But if you have $(7, 7)$ matrix where chances of overfitting are high you can make keep-prob lower.

3) *

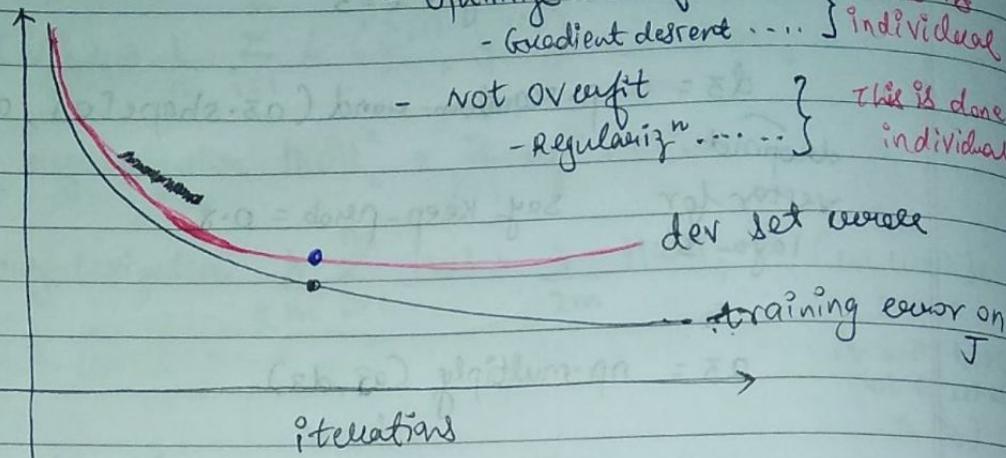
Data-Augmentation

4 \rightarrow 4 4 4

you can use some image & perform some transform on it.

4) Early Stopping

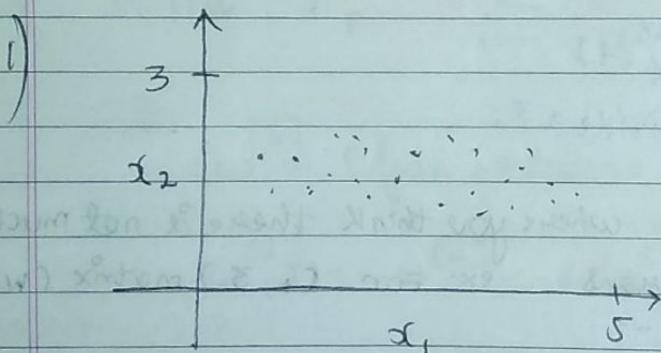
- optimize cost fn J } This is individual
- Gradient descent ...
- not overfit } This is done individual
- Regularizn' ...



You are stopping gradient descent in between.

By Early stopping, it mixes both optimiz'n & overfitting reduce methods.

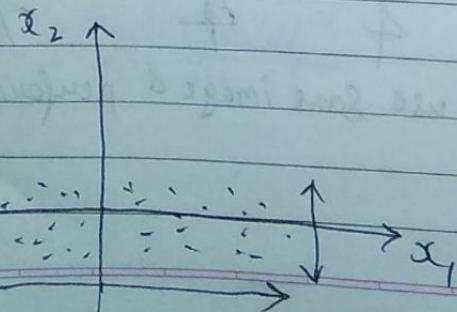
* To Speed up training. (NORMALIZATION)



Subtract mean:-

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x = x - \mu$$

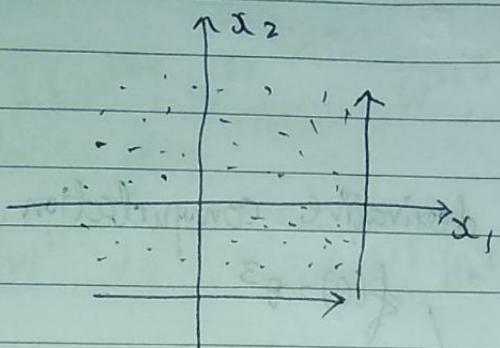


Normalize Variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot x^{(i)}$$

$$x' = \frac{x}{\sigma}$$

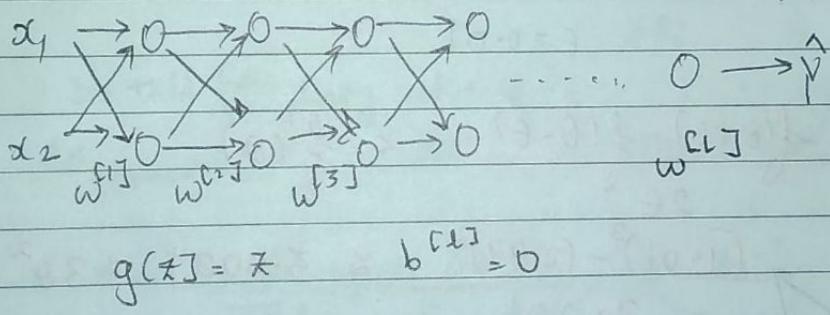
element wise scaling.



Now we have same variance of x_1 & x_2

If you use μ, σ^2 in training data, use same values in test set.

2). Vanishing / exploding gradients.



$$\hat{y} = w^{[L]} w^{[L-1]} \dots w^{[2]} w^{[1]} x$$

Assume each of this is $z^{[l]} = a^{[l]}$

$$w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1}$$

For deep layers { if $w^{[l]} > I$ (I is identity matrix)
 (Activations can explode exponentially)
 $w^{[l]} < I$
 (Activation can decrease exponentially)

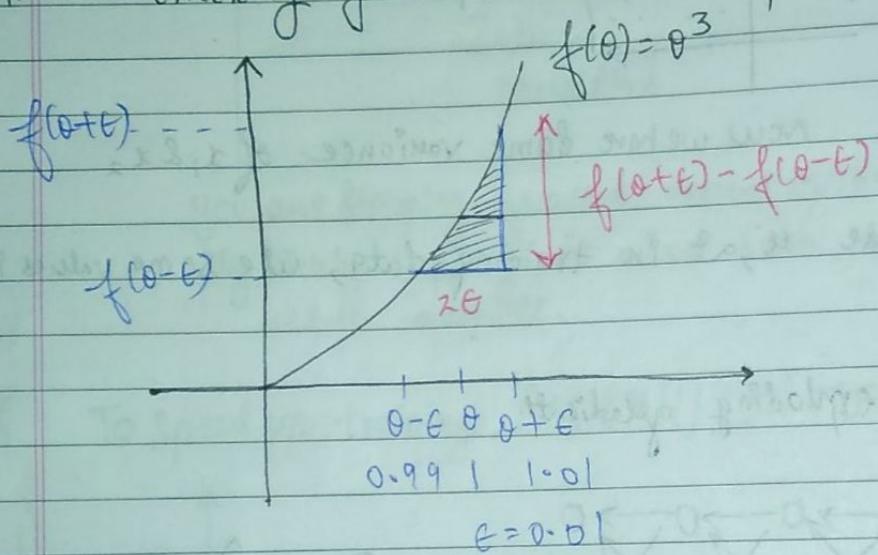
* weight initialization (To avoid exploding & vanishing)

$$w^{[l]} = np.random.randn(\text{shape}) * \sqrt{\frac{2}{n_{\text{in}} + 1}}$$

(For Relu)

$$\text{val}(w_i) = \frac{2}{n}$$

* Checking your derivative computation.



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx f'(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2 \times 0.01} \approx 3.0001 \approx 30^2 \approx 3$$

When you take

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx 3.0301$$

over 1 if 0.03 which is

so much

* Gradient checking.

Take $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$ & reshape into big vector θ .

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

similar for,

$d\theta^{[1]}, d\theta^{[2]}, \dots, d\theta^{[L]}$ into big vector $d(\theta)$

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each i :

$$\begin{aligned} d\theta_{approx}[i] &= \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon} \\ &\approx d\theta[i] = \frac{\delta J}{\delta \theta_i} \end{aligned}$$

Do this for every value of i ,

$d\theta_{approx} \approx d\theta$ (check whether they are equal?)

Check $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \leq 10^{-7}$ - great!

Take $\epsilon = 10^{-7}$ & you get

- * Don't use in training - only to debug
 - If algorithm fails, look at components to try to identify bug.
 - Remember regularization
 - Doesn't work with dropout

* Optimization algorithms *

DATE _____
PAGE NO. _____

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} \end{bmatrix}^T$$

$\underbrace{\quad\quad\quad}_{c_{n_x, m}} \quad \underbrace{\quad\quad\quad}_{x^{(1)} \dots x^{(1000)}} \quad \dots \quad x^{(m)}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} \end{bmatrix}^T$$

$\underbrace{\quad\quad\quad}_{(1, m)} \quad \underbrace{\quad\quad\quad}_{y^{(1)} \dots y^{(1000)}} \quad \dots \quad y^{(m)}$

This is if you have $m=5,000,000$

5,000 mini-batches of 1000 each

mini-batch $t: X^{\{t\}}, Y^{\{t\}}$

* How mini-batch works.

for $t=1, \dots, 5000$

Forward prop on $X^{\{t\}}$

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$A^{[L]} = g^{[L]}(Z^{[L]})$$

$$\text{Compute cost } J = \frac{1}{1000} \sum_{i=1}^d L(y^{(i)}, A^{[L]}) + \frac{\lambda}{2 \cdot 1000} \sum_{j=1}^d \|W_j\|_F^2$$

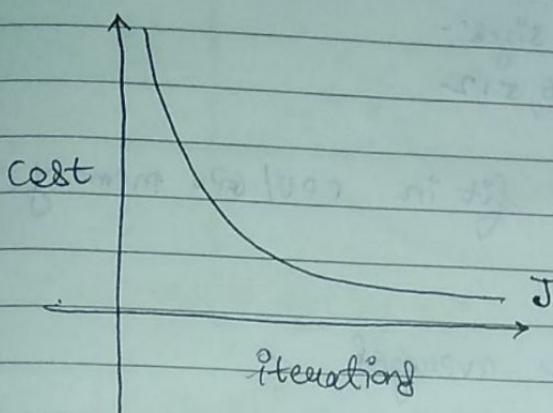
Back prop to get gradients w.r.t $J^{\{t\}}$ (using $(X^{\{t\}}, Y^{\{t\}})$)

$$W^{[L]} := W^{[L]} - \lambda dW^{[L]}$$

$$b^{[L]} := b^{[L]} - \lambda db^{[L]}$$

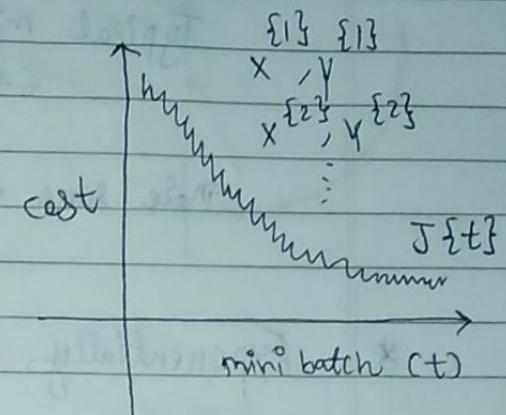
This is called "1 epoch"

For Batch



single pass through training set.

mini-batch



Plot $J^{(t)}$ computed using
 $x^{(t)}, y^{(t)}$

* Choosing your mini-batch size

If mini-batch size = m : Batch gradient descent $(x^{(1)}, y^{(1)}) = (x, y)$

If mini-batch size = 1 : Stochastic gradient descent

$$(x^{(1)}, y^{(1)}) = (x^{(1)}, y^{(1)})$$

so on.

Every example
is its own mini-
batch.

In practice: mini-batch size b/w 1 & m

Stochastic gradient
descent



lose spreading
from vectorization.

In betn

mini-batch

size (middle
not much
big or small)

Fastest
learning

Batch gradient descent

(mini-batch size = m)



Too much time
per iteration.

If small training set: use batch gradient descent
($m < 2000$)

Typical mini-batch sizes:-
64, 128, 256, 512

Make sure mini-batch fit in CPU/GPU memory.

* Exponentially weighted averages

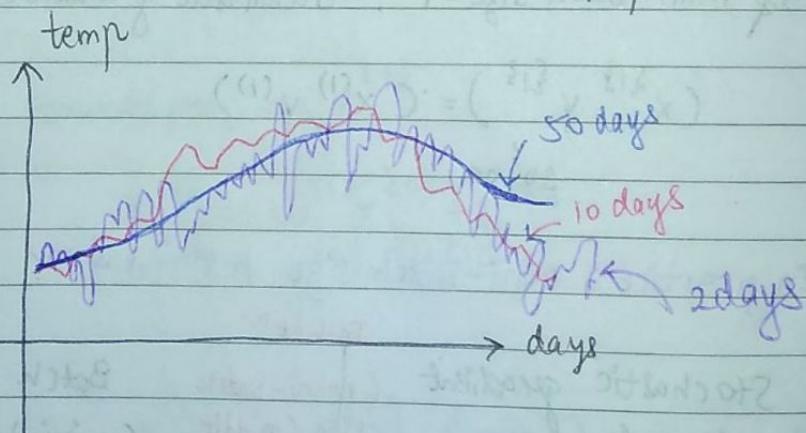
$$v_t = \beta v_{t-1} + (1-\beta) \theta_t$$

$\beta = 0.9$: ≈ 10 days temp.

$\beta = 0.98$: ≈ 50 days

$\beta = 0.5$: ≈ 2 days

v_t approximately averages over $\approx \frac{1}{1-\beta}$ days temp.



$$v_{100} = 0.9 v_{99} + 0.1 \theta_{100}$$

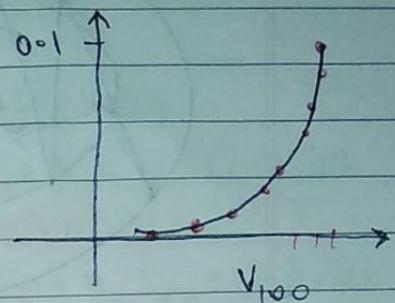
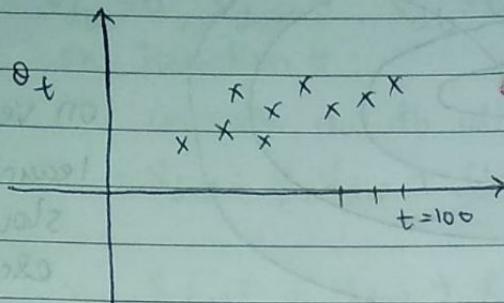
$$v_{99} = 0.9 v_{98} + 0.1 \theta_{99}$$

$$v_{98} = 0.9 v_{97} + 0.1 \theta_{98}$$

$$v_{100} = 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 (v_{98}))$$

$$= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 (0.9 v_{97} + 0.9 \theta_{98}) + \dots)$$

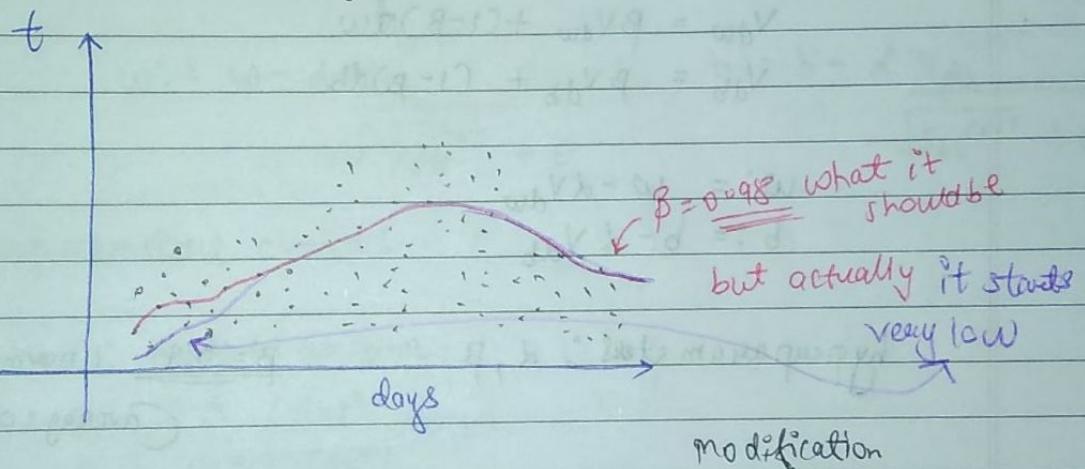
$$V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \times \theta_{99} + 0.1 \times (0.9)^2 \times \theta_{98} + 0.1 \times (0.9)^3 \times \theta_{97} + \dots$$



This is for $\beta = 0.9$

$$(0.9)^{10} \approx 0.35 \approx \frac{1}{e}$$

* Bias correction in weighted averages.



$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_0 = 0$$

$$V_1 = 0.98 V_0 + 0.02 \theta_1$$

$$V_2 = 0.02 \theta_2 + 0.98 V_1$$

$$= 0.98(0.02 \theta_1) + 0.02 \theta_2$$

$$V_2 = 0.0196 \theta_1 + 0.02 \theta_2$$

modification

$$V_t = \frac{V_{t-1}}{1-\beta^t}$$

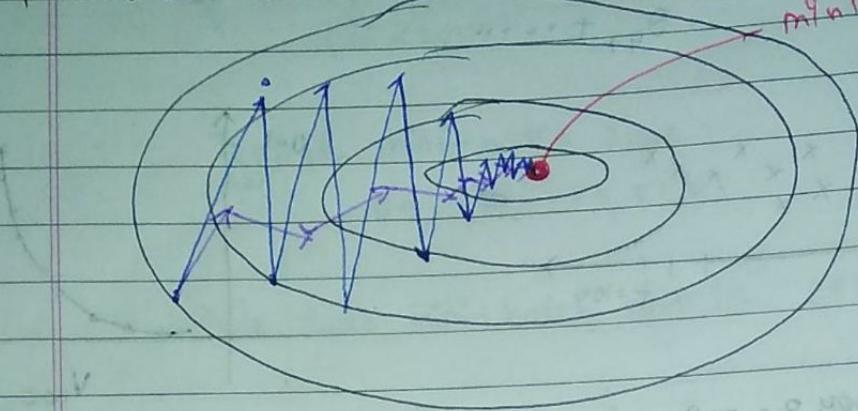
$$t=2 \therefore 1-\beta^t = 1-0.98^2 \\ = 0.0396$$

$$\therefore V_2 = \frac{V_1}{0.0396}$$

$$V_2 = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

So when t is very large, β^t doesn't make any diff. so it helps for initial temperatures.

* Gradient-descent with momentum



on vertical axis,
learning should be
slower to reduce
oscillations.

on horizontal axis,
faster learning.

Momentum :-

on iteration t :

compute dw, db on current mini-batch

$$v_{dw} = \beta v_{dw} + (1-\beta)dw$$

$$v_{db} = \beta v_{db} + (1-\beta)db$$

$$w := w - \alpha v_{dw}$$

$$b := b - \alpha v_{db}$$

hyperparameters: α, β

$\beta = 0.9$ (normally used)

(average over last 10
gradients)

* RMSprop

on iteration t :

compute dw, db on current mini-batch

element wise

$$s_{dw} = \beta s_{dw} + (1-\beta)dw^2 \leftarrow \text{small}$$

$$s_{db} = \beta s_{db} + (1-\beta)db^2 \leftarrow \text{large}$$

Ex:- we want to take small & w take large

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

$\sqrt{s_{dw} + \epsilon}$

small so large to speed up in w direction

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

$\sqrt{s_{db} + \epsilon}$

large & so slow down in b direction

for numerical stability $\epsilon = 10^{-8}$

* Adam Optimization algorithm. :

$$v_{dw} = 0, S_{dw} = 0, v_{db} = 0, S_{db} = 0$$

on iteration t:

compute dw, db using current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw, v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

"momentum" β_1

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2, S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

"RMS prop" β_2

$$w := w - \frac{v_{dw}}{\sqrt{S_{dw}} + \epsilon}, b := b - \frac{v_{db}}{\sqrt{S_{db}} + \epsilon}$$

$$v_{dw} = \frac{v_{dw}}{\sqrt{S_{dw}} + \epsilon}, v_{db} = \frac{v_{db}}{\sqrt{S_{db}} + \epsilon}$$

$$w := w - \frac{v_{dw}}{\sqrt{S_{dw}} + \epsilon}, b := b - \frac{v_{db}}{\sqrt{S_{db}} + \epsilon}$$

hyperparameter choice:

α : needs to be fine

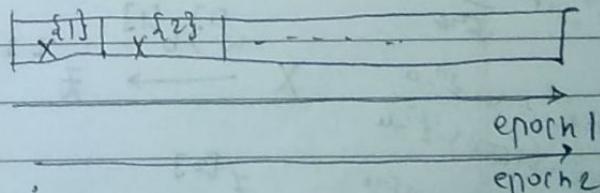
β_1 : 0.9 (dw)

β_2 : 0.999 (dw^2)

ϵ : 10^{-8}

* Learning rate decay.

1 epoch = 1 pass through data



$$\alpha = \frac{1}{1 + \text{decay rate} * \text{epoch_num}}$$

$$\text{let } \alpha_0 = 0.2$$

$$\text{decay rate} = 1$$

Epoch	α
1	0.1
2	0.67
3	0.5
⋮	⋮

* Implementing batch norm

Given some intermediate values in NN

$x^{(1)}, \dots, x^{(m)}$

$\bar{x}^{(i)}$

$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

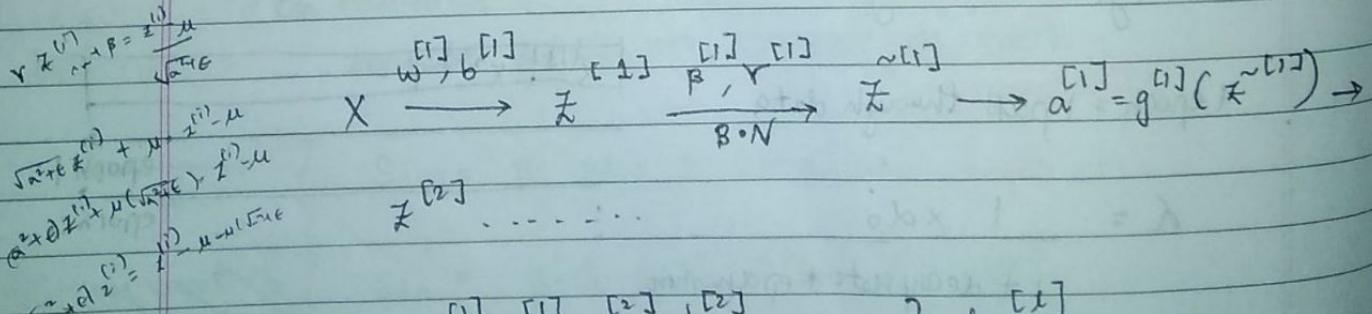
$$\sigma^2 = \frac{1}{m} \sum_i (x_i - \mu)^2$$

$$\begin{matrix} x^{(i)} \\ \text{norm} \end{matrix} = \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$x^{(i)} = \gamma \frac{x^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta \quad \begin{matrix} \text{learnable parameters} \\ \text{of model.} \end{matrix}$$

$$\text{If } \gamma = \sqrt{\sigma^2 + \epsilon} \text{ then } \underline{x^{(i)}} = \underline{x^{(i)}}$$

* Adding Batch norm to a network.



Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, \gamma^{[1]}, \beta^{[1]}, \gamma^{[2]}, \beta^{[2]}, \dots$ { $d\beta^{[1]}$
 $\beta^{[1]} = \beta^{[1]} - d\beta^{[1]}$

Multi-class classification.

* Softmax function

$$\tilde{z}^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \quad (4, 1)$$

Softmax Activation function:

$$\rightarrow t = e^{\tilde{z}^{[L]}} \quad (4, 1)$$

$$\rightarrow a^{[L]} = \frac{e^{\tilde{z}^{[L]}}}{\sum_{j=1}^4 t_j} = a, \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

(4, 1)

$$\tilde{z}^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$a^{[L]} = g^{[L]}(\tilde{z}^{[L]}) \\ = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix}$$

Softmax

$$= \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

$$\text{hard max} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad (\because 5 \text{ was biggest it gave 1 to it})$$