

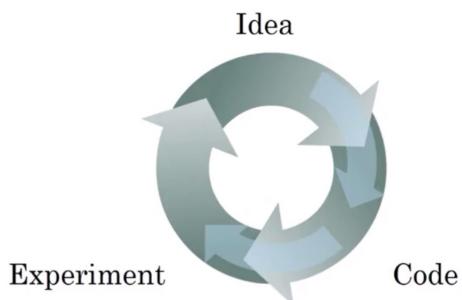
# Course 2: Improving Deep Neural Networks:

## Hyperparameter tuning, Regularisation and Optimisation

### Week 1: Practical aspects of Deep Learning

Applied ML is a highly iterative process

- # layers
- # hidden units
- learning rates
- activation functions
- ...



### Train/dev/test sets

Data

--	--	--

training  
set

hold-out set  
OR cross-validation set  
OR development set  
→ "dev set"

Previous Era:

70% - 30% split or 60% - 20% - 20%

Big data era:

% of dev set and test set can be much smaller.

98% - 1% - 1% (for large amount of data)

## Mismatched train/test distributions

e.g. a model to find picture of cars among user-uploaded picture.

training set:

cat pictures from  
web pages

↓  
nice, clear pic

Dev/test set:

cat pictures from users  
using your app

↓  
blur, casual pic

→  
distribution of  
data may be different

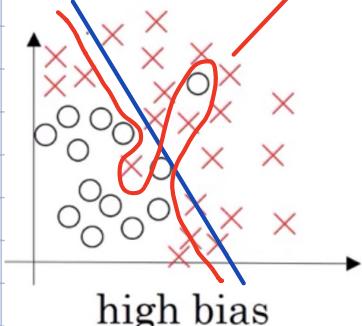
↳ make sure the dev/test set come from the same distribution

↳ we will be using the dev set to evaluate a lot of different models

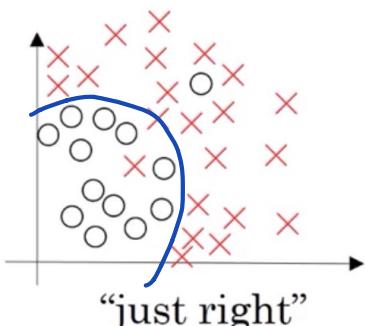
It might be ok to not have a test set (only dev set)

## Bias / Variance

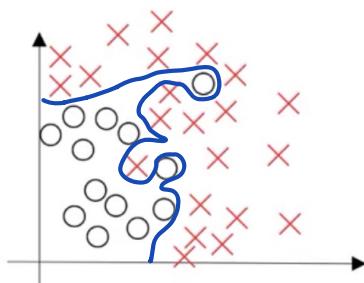
high bias  
and high variance



high bias



"just right"



high variance

underfitting

overfitting

In high dimension problem, can't visualize decision boundary



$y=1$

$y=0$

training set error:

1%

15%

15%

0.5%

dev set error:

11%

16%

30%

1%

high  
variance

(can't generalise  
well)

high  
bias

↓  
underfit

high  
bias

AND  
high  
variance

low  
bias  
low  
variance

Assume Human can achieve ~ 0% err

→ over optimal (Bayes) error ~ 0% err

→ really blur image  $\Rightarrow$  Bayes error will be higher

### Basic Recipe for ML

After trained an initial model,

first ask: does the algorithm have high bias?  $\rightarrow$  training set performance

$\hookrightarrow$  if yes  $\rightarrow$  bigger network / different architecture

(keep trying until no longer high bias)

i.e. at least fit the training set well

then ask: do we have a variance problem?  $\rightarrow$  dev set performance

$\hookrightarrow$  if yes  $\rightarrow$  more data / regularisation / architecture

Bias Variance trade-off.



in modern era, getting a bigger network almost always just reduce bias w/o hurting variance. Getting more data pretty much always reduces variance w/o hurting bias.

With regularisation, training a bigger network almost never hurts.

### Regularisation in neural network

If we find the network is over fitting the data, one of the first thing to try is probably regularisation. (another reliable way is to get more data)

Let's develop this idea using logistic regression.

recall we want to minimize the cost function  $J$ :

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

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

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

this is L2 regularisation

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

L1 regularisation

If we use L1 regularisation,  $w$  will end up being sparse.

$\hookrightarrow$  Lasso

$\lambda$  = regularisation parameter

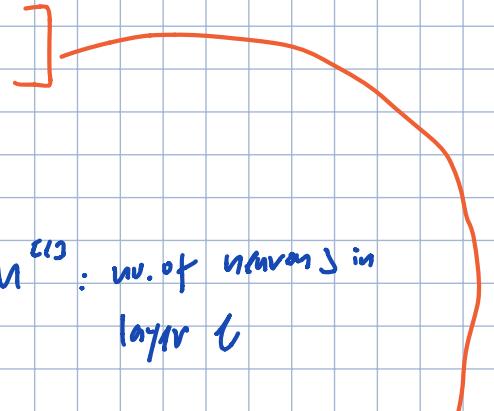
## How about neural network?

In a NN, we have a cost function that's a function of all our parameters:

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)})$$

$L$ : no. of layers

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



where the squared norm of a matrix

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_j (w_{ij}^{(l)})^2$$

$n^{(l)}$ : no. of neurons in layer  $l$

"Frobenius norm"

$$w: (n^{(0)}, n^{(L-1)})$$

↳ sum of square of elements of a matrix

How to compute gradient descent with this?

$$\rightarrow \text{Previously, we would complete } dw^{(l)} = (\text{back prop}) + \frac{\lambda}{m} w^{(l)}$$

where back prop gives us  $\frac{\partial J}{\partial w^{(l)}}$  ↪ still holds with regularisation

↳ then update

$$w^{(l)} := w^{(l)} - \alpha dw^{(l)}$$

w/o regularisation

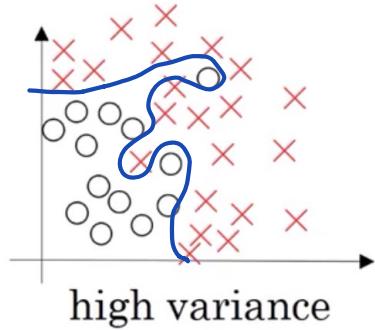
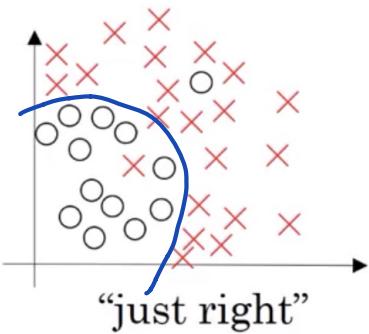
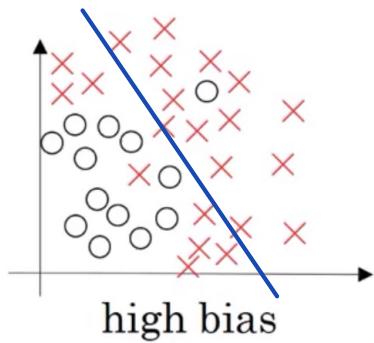
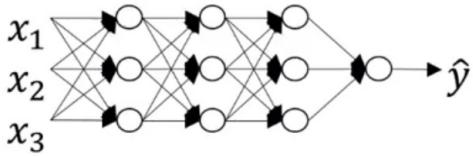
$$\text{L2 is also called weight decay: } w^{(l)} := w^{(l)} - \alpha [(\text{from back prop}) + \frac{\lambda}{m} w^{(l)}]$$

$$= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{from back prop})$$

$$= (1 - \frac{\alpha \lambda}{m}) w^{(l)} - \alpha (\text{from back prop})$$

with  
regularisation

## Why Regularisation reduces overfitting?



Say we have a neural network that is currently overfitting -

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

↖  
penalizes the weight matrices  
from being too large

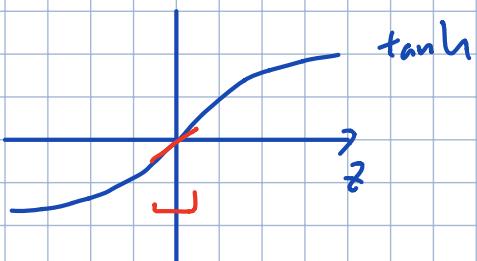
if  $\lambda$  is very large  $\rightarrow W^{[L]} \approx 0$

$\hookrightarrow$  a lot of neurons will be zero (sort of)

$\hookrightarrow$  become a simplified network

$\hookrightarrow$  less variance

Another intuition:



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

if  $z$  is small  $\rightarrow$  linear region of tanh is used

if  $\lambda$  is large, then  $W^{[l]}$  is small

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

$\hookrightarrow z$  is also small

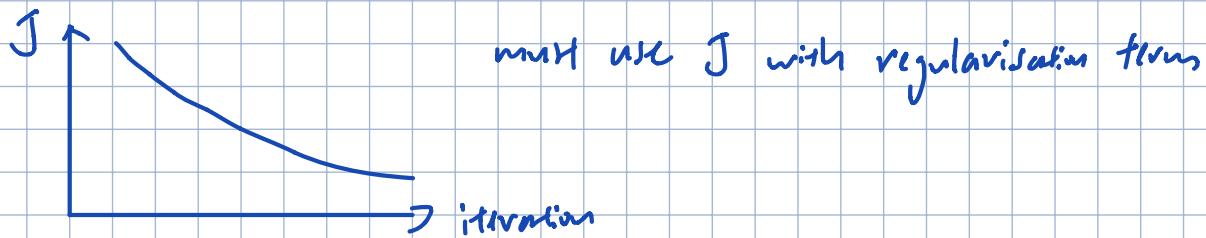
$\hookrightarrow g(z)$  is roughly linear

$\hookrightarrow$  every layer will be roughly linear

$\hookrightarrow$  whole network is linear

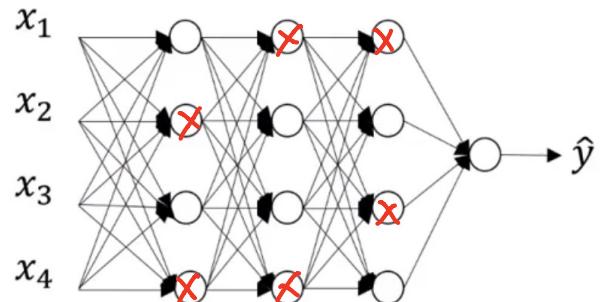
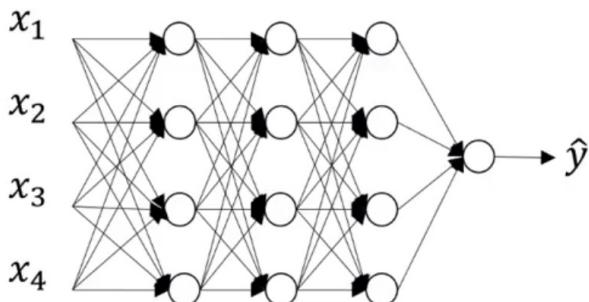
$\hookrightarrow$  less variance

when implementing gradient descent, can plot:

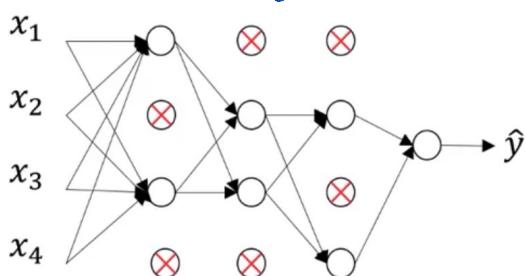


## Dropout Regularisation

Let's say we train a neural network and there's overfitting.



remove all outgoing also



With dropout, we go through each of the layers of the network and set some probability of eliminating a node in the network

e.g. for each node, we toss a coin, have a 0.5 chance of keeping each node and 0.5 chance of removing the node.

↳ actually works as regularisation, as the network is simpler.

There are a few ways of implementing dropout, the most common way is called "Inverted dropout"

Let's say we want to illustrate this with layer  $L=3$

we set a vector "d"

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

dropout  
vector  
for layer 3

each element: 0.8 chance = 1  
0.2 chance = 0

a number e.g. 0.2  
probability a given  
hidden unit will be  
kept

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

activation  
of layer 3

$a_3' = \text{keep-prob} \leftarrow \text{we need to scale-up}$

↳ 50 neurons in 3rd hidden layer

↳ 10 neurons shut off

$$z^{(4)} = W^{(4)} \cdot a^{(3)} + b^{(4)}$$

reduced by 20%

expected value reduced.

need scale-up

on various pass (e.g. gradient descent), we want to shrink out different neurons.

⇒ At test time, making prediction:

Given some  $a^{[0]} = X$

NO Dropout in prediction

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

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

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

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



⇒ if randomly dropout in test time → add noise only

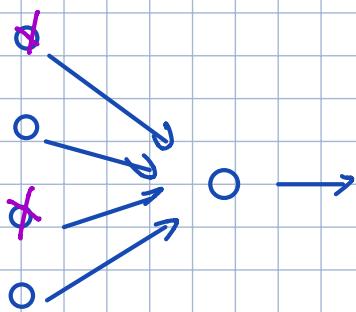


### Understanding Dropout

Apart from resulting a smaller network, which has regularization effect, here is another intuition

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

Let's look at the perspective of a single unit.



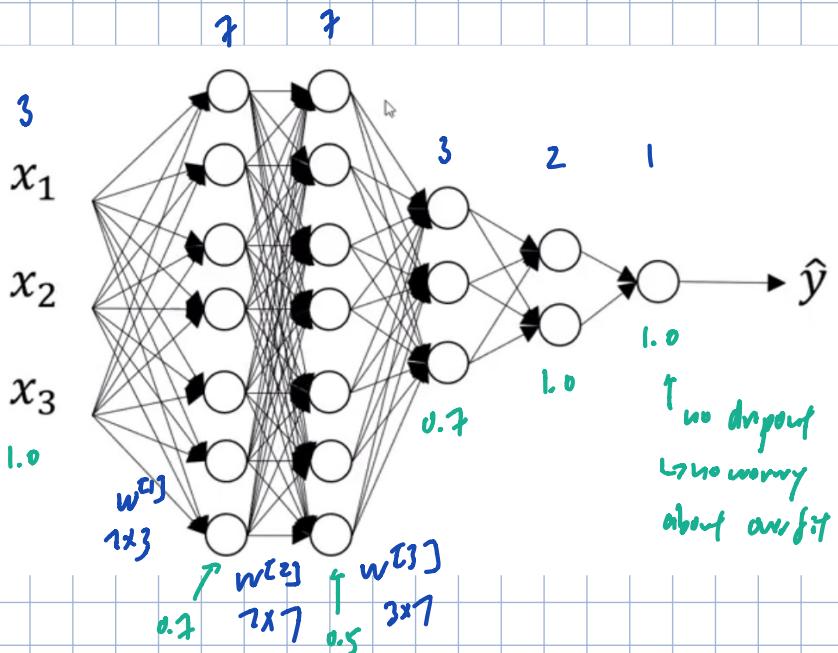
for dropout regularisation, input can get randomly eliminated

↳ this neuron can't rely on any one feature

↳ can't put a lot of weight on any one input

↳ motivated to spread out the weight

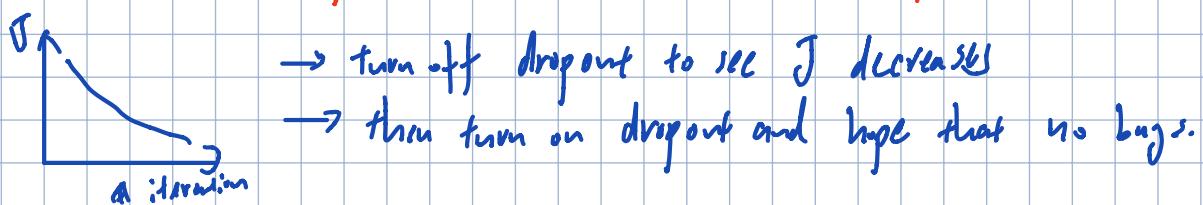
↳ shrinking the squared norm of weights  
similar to L2



no. of neurons in previous layer  
 = no. of col of w  
 no. of neurons in current layer  
 = no. of row of w

Not all layer uses dropout, because  $w^{(2)}$  is the largest  $\Rightarrow$  maybe apply dropout on  $w^{(2)}$   
 so for layer 2, keep-prob = 0.5

Downside: cost function  $J$  is not well-defined with dropout.  
 $\hookrightarrow$  on every iteration, we randomly kill off a bunch of nodes.

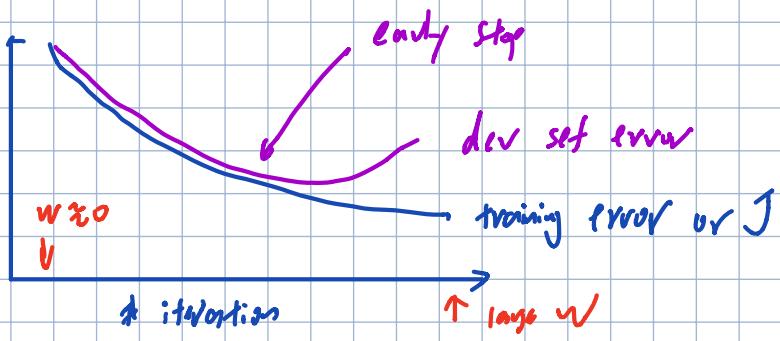


### Other regularisation method

- ① Data augmentation  $\rightarrow$  create full training examples  
 $\hookrightarrow$  more data  $\rightarrow$  sort of regularise  $\rightarrow$  reduce overfitting

### ② Early Stopping

random init  
 $\hookrightarrow$   $w$  to small value



Early stop gives w/ with medium size values

↳ similar to L2 regularisation

Down side:

we suppose to optimise cost function  $J \rightarrow$  gradient descent, ...

↳ also not overfitting  $\rightarrow$  Regularisation, ...

→ 2 task separately, use different w/  $\rightarrow$  orthogonalisation

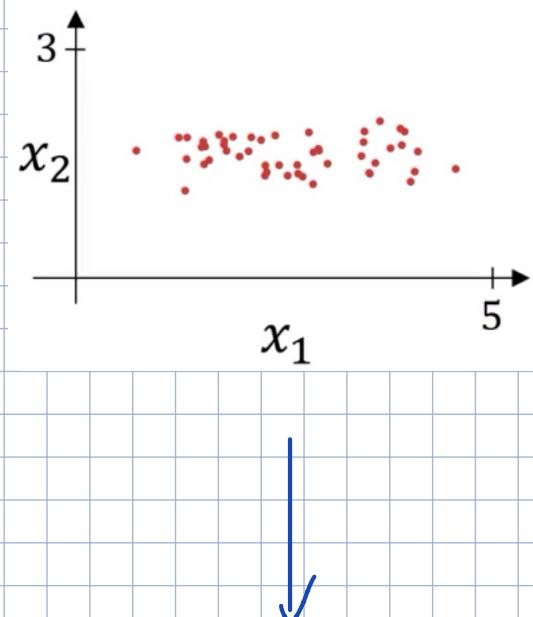
Early stop couples these 2 tasks

### Setting up optimisation problem

#### Normalising Input

This is to speed up training

e.g. training set with 2 input features



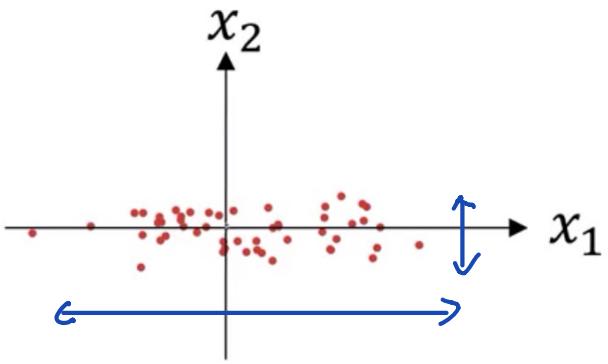
$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Normalize inputs involves 2 steps:

0 Subtract mean:

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

$$x := x - \bar{x}$$



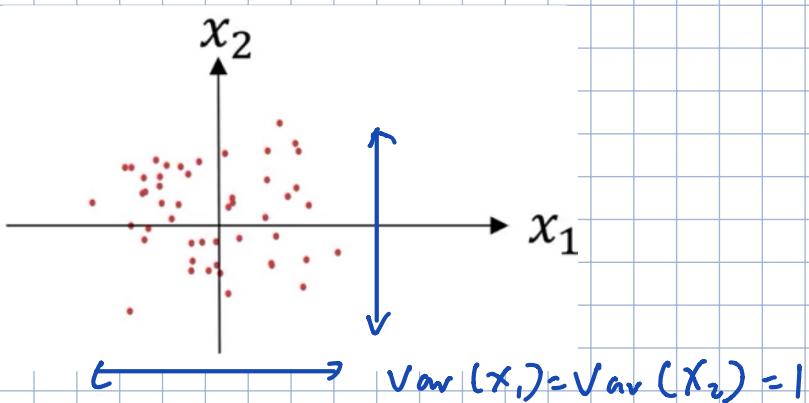
③ normalize the variance

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

$\sigma^2$  element wise

*we already  
subtracted mean  
in D*

$$\frac{x - \bar{x}}{\sigma}$$



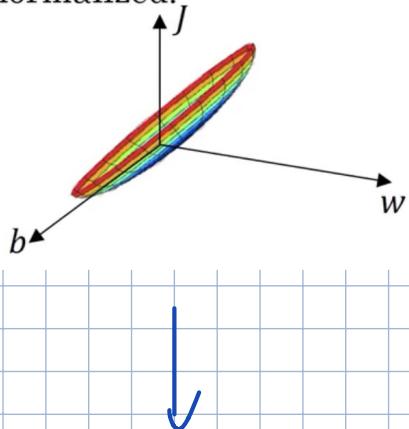
\* If we use this to scale training data

$\Rightarrow$  use the same  $\bar{x}$  and  $\sigma$  to scale test data!

why normalize inputs?

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

Unnormalized:

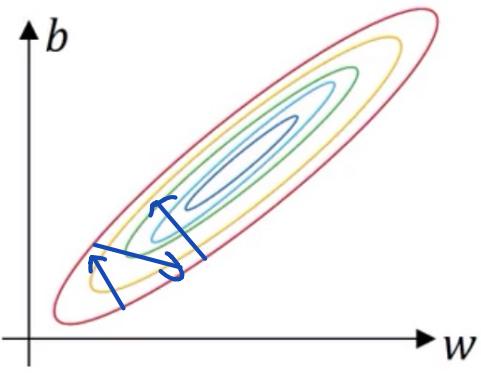


if  $x_1 : 1, \dots, 1000$

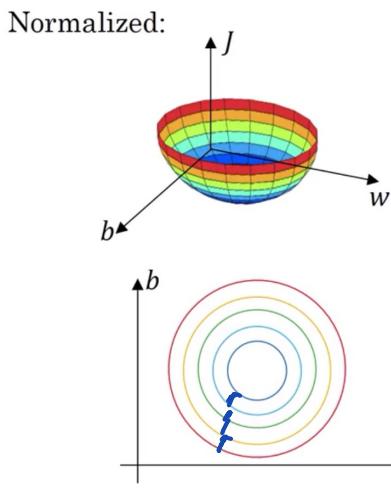
$x_2 : 0, \dots, 1$

$\hookrightarrow w_1$  and  $w_2$  will have very  
different value

$\hookrightarrow$  elongated bowl



need small learning  
rate of gradient  
descent



can use larger learning rate

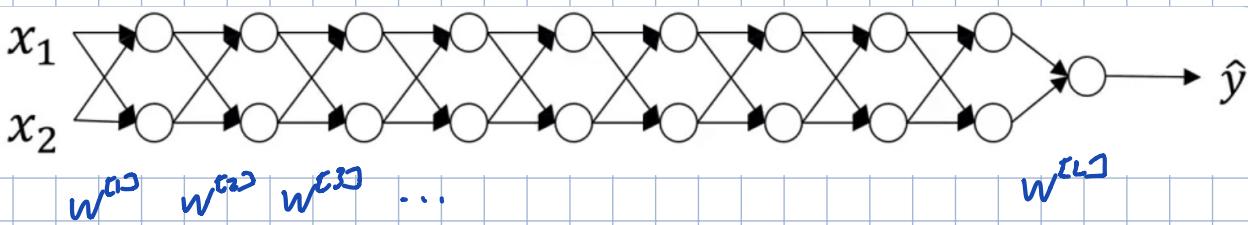
### Vanishing / Exploding gradients

One of the problems when training NN, especially very deep NN is about vanishing and exploding gradients.

i.e. the derivatives / slopes can get very big or small

↳ use careful choices of the random weight initialisation to significantly reduce this problem

Let's say we train a very deep NN like this:



Say for simplification, we use linear activation

$$g(z) = z$$

and ignore  $b$ , i.e.  $b^{(L)} = 0$

$$\text{we can get } \hat{y} = w^{(L)} w^{(L-1)} \dots w^{(2)} w^{(1)} x$$

$z^{(1)} = w^{(1)} x$

Say each  $w^{[l]}$  is a bit larger than identity

$$w^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$

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

$\hookrightarrow$  if  $L$  is large, then  $\hat{y}$  will be large

If  $w^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$

$$\hookrightarrow \hat{y} = w^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x$$

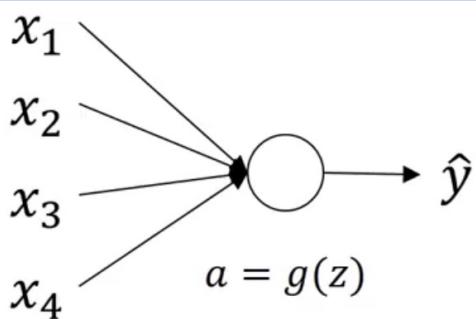
$\hookrightarrow$  if  $L$  is large, then  $\hat{y}$  will be small

$\Rightarrow$  can be inferred to derivatives

if gradient is very small  $\rightarrow$  takes long for gradient descent

## Weight initialisation for deep networks

single neuron example



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b^0$$

larger  $n \longrightarrow$  smaller  $w_i$

$\because z$  is sum of  $w_i x_i$ , so if  $n$  is large  
we want each term be smaller

$$\text{we can set } \text{Var}(w_i) = \frac{1}{n}$$

where  $n$  is # of input feature into a neuron

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

if activation is ReLU, then we  $\frac{2}{n_{\text{in}}}$



So if the input features or activations are roughly mean 0 and standard variance  $\rightarrow$  results  $w$  not too large / smaller than 1.

Other variants:

if tanh is used

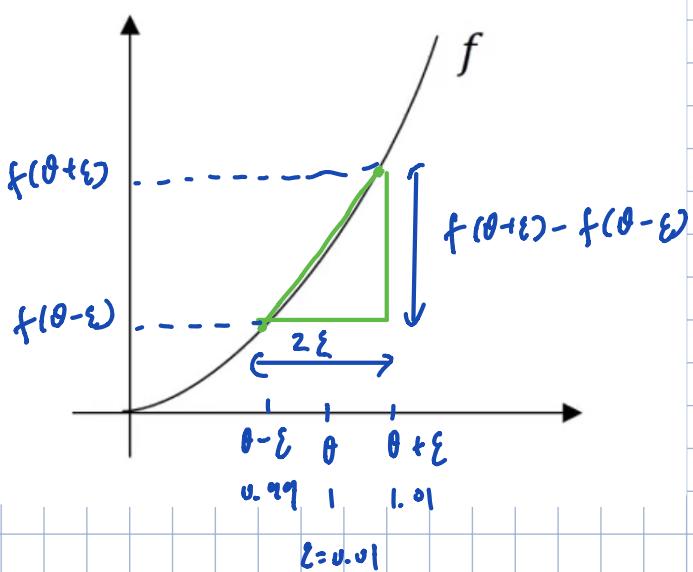
$\hookrightarrow \sqrt{\frac{1}{n^{(L-1)}}}$  is used (Xavier initialisation)

or  $\sqrt{\frac{2}{n^{(L-1)} + n^L}}$

### Numerical approximation of gradients

when implementing back propagation, there is a test called gradient checking, that can help checking the back prop is correct.

$\hookrightarrow$  first, need approximation of gradients



$$\text{say } f(\theta) = \theta^3$$

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

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.001

$$\text{formally: } f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \rightarrow \text{error } O(\epsilon^2)$$

$$\text{if } \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{\epsilon} \rightarrow \text{error } O(\epsilon)$$

## Gradient Checking

The NN will have some parameters:

take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  and reshape into a big vector  $\theta$

$$\hookrightarrow J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) \rightarrow J(\theta) \quad \text{concatenate}$$

then, take  $d\theta^{(1)}, d\theta^{(2)}, \dots, d\theta^{(L)}$  and reshape into a big vector  $d\theta$

( $\theta$  and  $d\theta$  has the same dimension)

→ Is  $d\theta$  the gradient/slope of  $J$ ?

Grad check:  $(\text{recall } J(\theta) = J(\theta_1, \theta_2, \dots))$

for each  $i$ : i.e. each component of  $\theta$

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i}$$

↳ end up with  $d\theta_{\text{approx}}$ , same dimension as  $d\theta$

check:  $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta\|_2}$  ↪ 0.2 norm: sum of squares of elements of the difference, then sqrt

$\frac{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}{\|d\theta\|_2}$  ↪ sum of squares of elements then sqrt

$$\approx \epsilon = 10^{-7}$$

$$\downarrow \text{if } \approx 10^{-7} \quad \text{good}$$

$$10^{-5} \quad \text{human ...}$$

$$10^{-3} \quad \text{wrong}$$

## Gradiant checking implementation Notes

- Do Not use in training → only to debug
  - i.e. don't calculate  $d\theta_{approx}[i]$  in training
  - Just get  $d\theta$  during back prop
- ↳ do not run AND in every iteration of gradient descent
- If algorithm fails grad check, look at components to try identify the bug
- Remember regularisation term, as it is part of  $J$ 
  - ↳  $d\theta = \text{grad of } J \text{ w.r.t. } \theta$
- Does NOT work with dropout
  - ↳ implement grad check w/o dropout
  - ↳ then turn off dropout
- Run at random init; then maybe again after some training
  - ↳ sometimes (rare), gradient descent is correct when  $W$  and  $b$  are close to 0, as  $w, b$  gets bigger, might not be correct.

## Week 2: optimisation algorithms

### Mini-batch gradient descent

Deep NN works well with large data set

↳ but train on large data is slow

### Batch vs mini-batch gradient descent

Vectorisation allows you to efficiently compute on  $m$  examples

$$X = [x^{(1)} \ x^{(2)} \ \dots \ x^{(1000)} \ | \ x^{(1001)} \ \dots \ x^{(2000)} \ | \ \dots \ | \ \dots \ x^{(m)}]$$

$(n_x, m)$        $x^{t, 1}$        $x^{t, 2}$        $\dots$        $x^{t, 1000}$        $x^{t, 1001}$        $\dots$        $x^{t, 2000}$        $\dots$        $x^{t, m}$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(1000)} \ | \ y^{(1001)} \ \dots \ y^{(2000)} \ | \ \dots \ | \ \dots \ y^{(m)}]$$

$(1, m)$        $y^{t, 1}$        $y^{t, 2}$        $\dots$        $y^{t, 1000}$        $y^{t, 1001}$        $\dots$        $y^{t, 2000}$        $\dots$        $y^{t, m}$

what if  $m = 5,000,000$ ?

traditionally, train  $5\text{ mil} \rightarrow 1$  step of gradient descent

but we can split to mini batches (say 1,000 examples each)

↳ can train on minibatch and progress through gradient descent

⇒ we have 5000 minibatches

mini batch no.  $t$ :  $X^{t, 1}, Y^{t, 1}$

$X^{(i)}$ :  $i^{\text{th}}$  training example

$Z^{t, 1}$ :  $l^{\text{th}}$  layer,  $z$  value

$X^{t, 1}$ :  $t^{\text{th}}$  mini batch train example

$X^{t, 1}$ :  $(n_x, 1000)$

$Y^{t, 1}$ :  $(1, 1000)$

## mini-batch gradient descent:

for  $t = 1, \dots, 5000$  {

Forward prop on  $X^{[t]}$

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

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

:

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

$$\text{compute } (1) J^{[t]} = \frac{1}{1000} \sum_{i=1}^C L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_C \|W^{[C]}\|_F^2$$

in this for loop  
implement 1 step  
of gradient descent  
using  $X^{[t]}, Y^{[t]}$

verbalized

for  $X^{[t]}, Y^{[t]}$

Backprop to compute gradient wrt  $J^{[t]}$

(using  $X^{[t]}, Y^{[t]}$ )

$$\text{update: } W^{[C]} = W^{[C]} - \alpha \nabla W^{[C]}$$

$$b^{[C]} = b^{[C]} - \alpha \nabla b^{[C]}$$

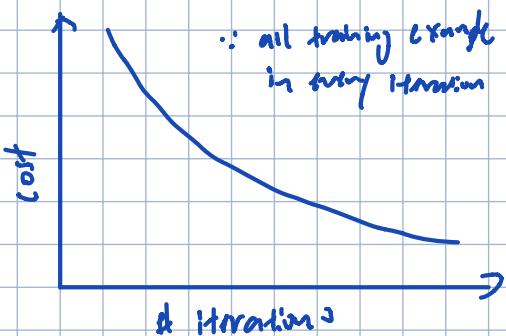
}

"1 epoch" of training

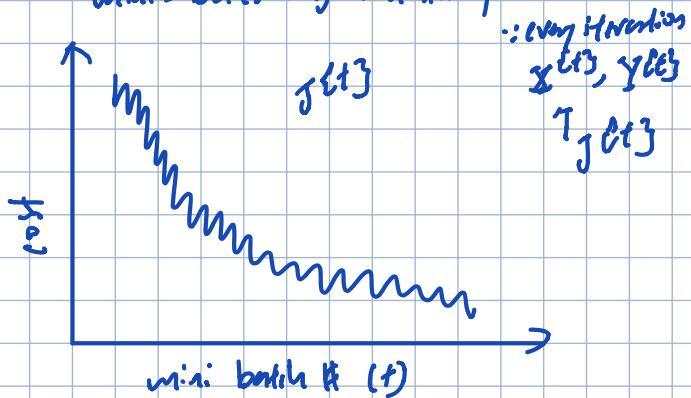
[repeat for multiple pass of the whole batch]

## Understanding mini-batch gradient descent

### Batch gradient descent



### mini-batch gradient descent



each mini-batch might easy/hard  
to train  $\rightarrow$  low/high  $J^{[t]}$

## Choosing the size of mini-batch

If mini-batch size =  $m$

↳ batch gradient descent  $\rightarrow$  just 1 mini-batch  $(\mathbf{X}^{(1)}, \mathbf{Y}^{(1)}) = (\mathbf{X}, \mathbf{Y})$

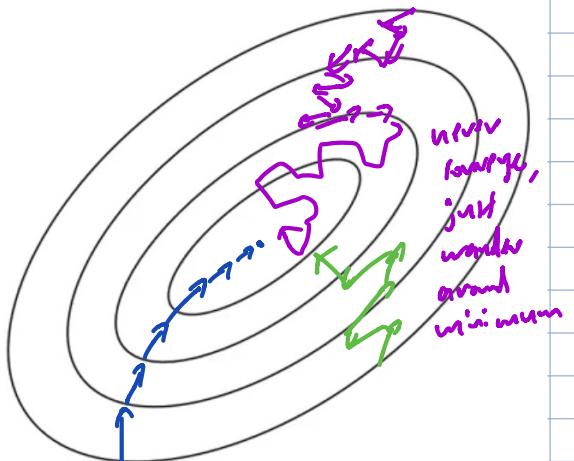
If mini-batch size = 1

↳ stochastic gradient descent

i.e. every example is its own mini-batch

$(\mathbf{X}^{(1)}, \mathbf{Y}^{(1)}) = (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}) \dots$

contour of  $J$



in-between: Faster Learning

↳ get a lot of iteration

↳ make progress w/o entire training set

In practice: mini-batch size will be somewhere in between 1 and  $m$

Batch gradient descent (mini-batch size =  $m$ )

↳ too long per iteration

Stochastic gradient descent (mini-batch size = 1)

↳ lose speed up from vectorisation

## Some guidelines

→ If we have a small training set, just use batch gradient descent  
(e.g.  $m < 2000$ )

→ Typical minibatch size: 64, 128, 256, 512 (power of 2, might run faster)

- make sure all of  $X^{[t]}$ ,  $y^{[t]}$  fits in CPU/GPU memory
- minibatch size is another hyper parameter

### Exponentially weighted averages

There are a few algorithms that are frequently used, and faster than gradient descent. Need to know exponentially (moving) weighted averages.

## Temperature in London

$$\theta_1 = 40^\circ\text{F} \quad 4^\circ\text{C} \quad 1 \text{ Jan}$$

$$\theta_2 = 49^\circ\text{F} \quad 9^\circ\text{C} \quad :$$

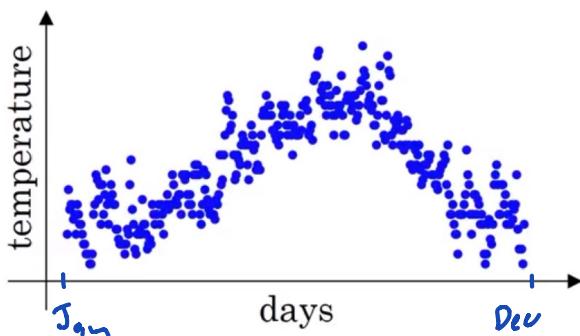
$$\theta_3 = 45^\circ\text{F} \quad :$$

:

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F} \quad :$$

:



### Moving average

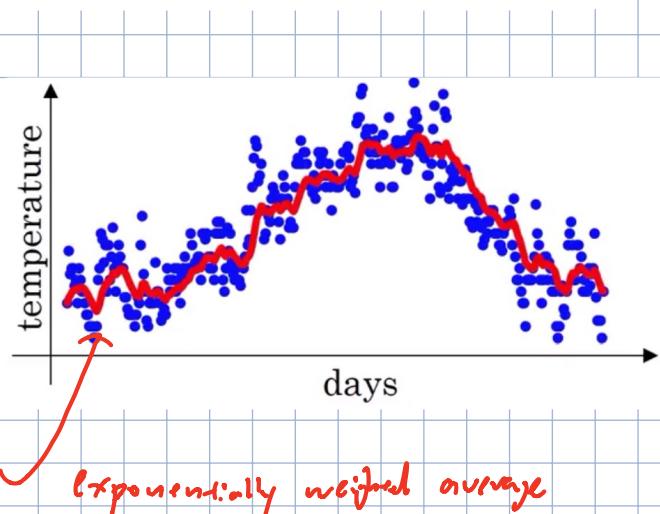
$$\text{First initial: } V_0 = 0$$

$$\hookrightarrow V_1 = 0.9V_0 + 0.1\theta_1$$

$$V_2 = 0.9V_1 + 0.1\theta_2$$

$$V_3 = 0.9V_2 + 0.1\theta_3$$

$$\rightarrow V_t = 0.9V_{t-1} + 0.1\theta_t$$



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

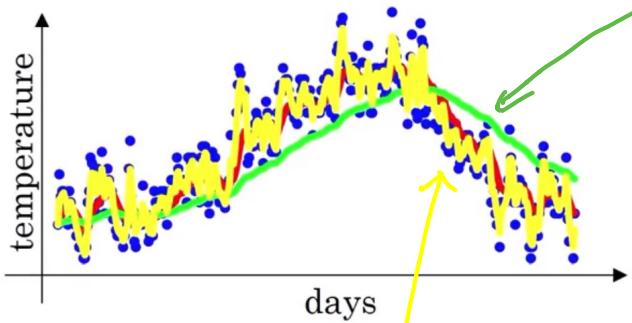
We can think  $V_t$  as approximately average over

$$\approx \frac{1}{1-\beta} \text{ days of temperature}$$

$\Rightarrow$  if  $\beta = 0.9 \Rightarrow \approx$  average over 10 days of temperature

$$\text{if } \beta = 0.98 \Rightarrow \frac{1}{1-0.98} \approx 50$$

$\hookrightarrow$  average over last 50 days of temperature



values like 0.98 gives a lot of weight on previous value  
 $\hookrightarrow$  more latency  
 $\hookrightarrow$  adapts more slowly

$$\beta = 0.5 \Rightarrow \frac{1}{1-0.5} = 2$$

$\hookrightarrow$  last 2 days average

$\hookrightarrow$  more susceptible to noise

$\hookrightarrow$  but adapt quickly

$\beta$  is another hyperparameter.

### Understanding exponentially weighted averages

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

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

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

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

:

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

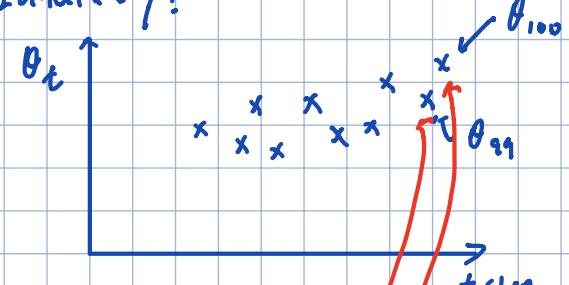
$$= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 V_{98})$$

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

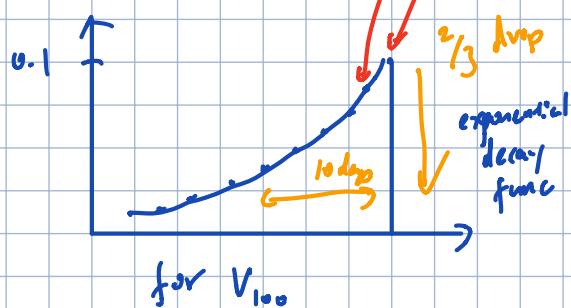
$$= \underline{0.1 \theta_{100}} + \underline{0.1 \times 0.9} \theta_{99} + \underline{0.1 \times (0.9)^2} \theta_{98} + \underline{0.1 \times (0.9)^3} \theta_{97} + \dots$$

All coeff add up to (or close to) 1

Intuitively:



so it is taking element wise product between these 2 function and sum it up



$$\therefore 0.1 \times (0.9)^n$$

$$0.9^{10} \approx 0.35 \approx 1/3$$

more generally:

$$(1-\varepsilon)^{1/\varepsilon} \approx 1/e$$

≈ It takes about 10 days for the height of decay from decay to 1/3

If  $\beta = 0.98$ , then  $0.98^{50} \approx 1/e \Rightarrow$  takes about 50 days

### Implementation of exponentially weighted averages

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3$$

...

$$V_\theta = 0$$

$$V_\theta := \beta V + (1-\beta) \theta_1$$

$$V_\theta := \beta V + (1-\beta) \theta_2$$

...

equivalently

$$V_\theta = 0$$

Repeat {

Get next  $\theta_t$

$$V_\theta := \beta V_\theta + (1-\beta) \theta_t$$

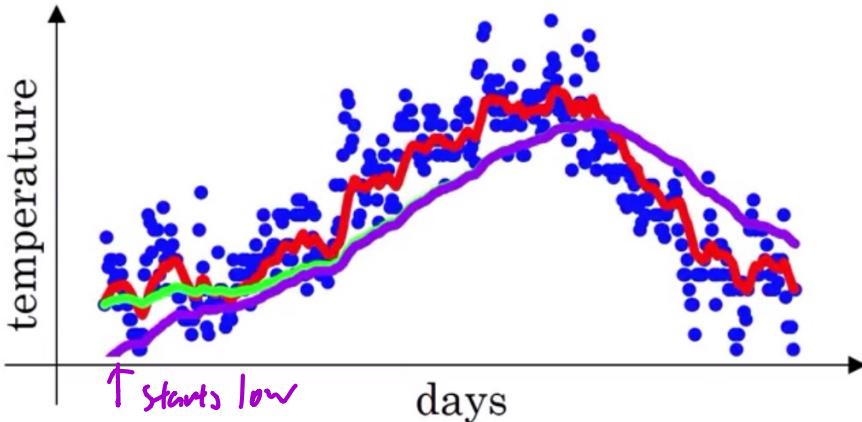
}

Advantage: just 1 variable to track

it moving window, might need to store the last 10 or 50 values

## Bias correction in exponentially weighted averages

Can be used to compute the averages more accurately.



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

we will get the purple line, instead of the green line (for  $\beta = 0.98$ )

When implementing a moving average, we initialise it

$$V_0 = 0 \quad \because V_0 \approx 0$$

$$\text{then } V_1 = \cancel{0.98V_0} + 0.02\theta_1 \leftarrow \text{low value at start}$$

Hence, to make this averaging more accurate, esp for the initial phase

$$\text{use } \frac{V_t}{1-\beta^t} \text{ instead of } V_t$$

$$\text{i.e. if } t=2, \quad 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\text{so estimate on day 2} = V_2 / 0.0396 = \frac{0.98V_1 + 0.02\theta_2}{0.0396}$$

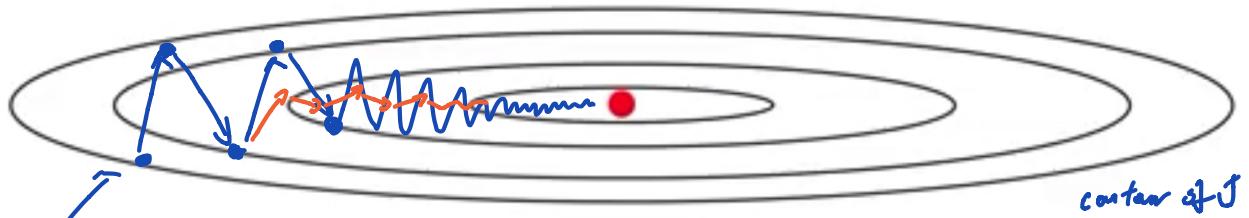
$$= \frac{0.98 \cdot 0.02\theta_1 + 0.02\theta_2}{0.0396}$$

$$\text{as } t \rightarrow \infty, \beta^t \rightarrow 0$$

## Gradient Descent with Momentum

Almost always works faster than the standard gradient descent algorithms

Basic idea: compute an exponentially weighted average of our gradients, then use that gradient to update our weights instead



standard gradient  
descent: oscillation  
prevents us using  
a large learning rate  
we might want  
↑ learn slower  
← learn faster

with momentum!

on iteration  $t$ :

compute the usual  $dW, db$  on current mini-batch

$$V_{dw} = \beta V_{dw} + (1 - \beta) dW \quad \begin{matrix} \text{trillion} \\ y \end{matrix} \text{ of velocity}$$

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

(iii) bowl-shape  $J$

provides acceleration

$$W := W - \alpha V_{dw}$$

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

} smoothing the raw gradients  
that is perp to descent  
direction

### Implementation details

$V_{dw} = 0$  same dimension as  $dW$ ,  $V_{db} = 0$  same dimension as  $db$  and  $b$

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

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

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

$$W = W - \alpha v_{dw}, \quad b = b - \alpha v_{db}$$

no bias correction  
needed in most  
case

Hyperparameters:  $\alpha, \beta$        $\beta = 0.9$

$\hookrightarrow \alpha > 0.9$   
Leverage over last 10 iterations gradient

## RMS prop

Also can be used to speed up gradient descent.

On iteration  $t$ :

compute the usual  $dW, db$  on current mini-batch

$$S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2$$

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

$$W := W - \frac{dW}{\sqrt{S_{dW}} + \epsilon}$$

$$b := b - \frac{db}{\sqrt{S_{db}} + \epsilon}$$

$\beta_2$  for RMS prop parameter

$\beta_1$  for momentum

make sure  
don't blow up (if  $S_{dW}/S_{db} \approx 0$ )

Very similar effect as momentum  $\rightarrow$  can use a larger learning rate  $\alpha$

## Adam Optimisation Algorithm (Adaptive Momentum Estimation)

Initialise  $V_{dW} = 0, S_{dW} = 0, V_{db} = 0, S_{db} = 0$

On iteration  $t$ :

momentum exponentially weighted average  $\rightarrow$  compute  $dW, db$  using current mini-batch

$$V_{dW} = \beta_1 V_{dW} + (1 - \beta_1) dW; V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

RMSprop  $\rightarrow S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2; S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$

$$\text{bias correction} \rightarrow V_{dW}^{\text{corrected}} = \frac{V_{dW}}{1 - \beta_1^t}; V_{db}^{\text{corrected}} = \frac{V_{db}}{1 - \beta_1^t}$$

$$S_{dW}^{\text{corrected}} = \frac{S_{dW}}{1 - \beta_2^t}; S_{db}^{\text{corrected}} = \frac{S_{db}}{1 - \beta_2^t}$$

$$W := W - \frac{V_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}}} + \epsilon} ; b := b - \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

## Hyperparameters choice of Adam:

$\alpha$ : needs to be tuned

$\beta_1$ : commonly use 0.9

( $d_w$ ) moving average

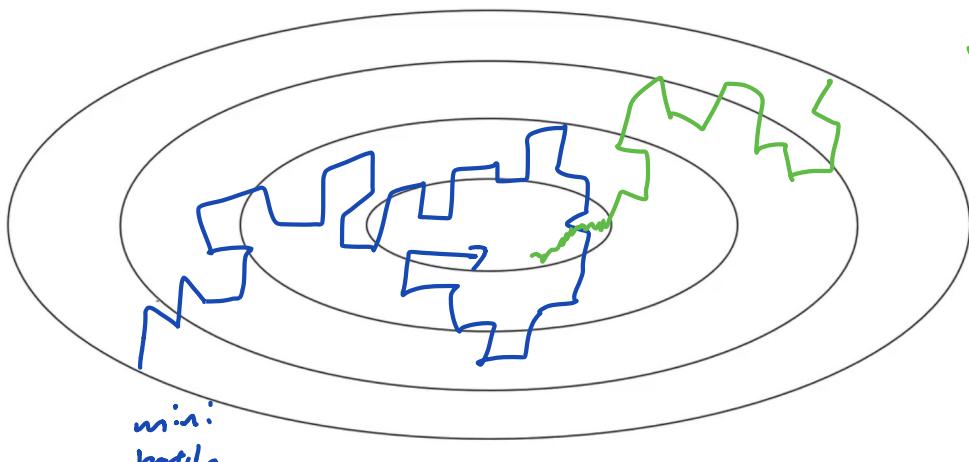
$\beta_2$ : 0.999

( $d_w^2$ ) moving average

$\epsilon$ :  $10^{-8}$

## Learning Rate Decay

Slowly reduce learning rate over time  $\rightarrow$  might speed up learning



slowly reduce  $\alpha$   
 $\hookrightarrow$  oscillate often  
near minimum

1 epoch = 1 pass through the data



$$\text{set } \alpha = \frac{\alpha_0}{1 + \text{decay\_rate} \times \text{epoch\_num}}$$

$\alpha_0$  initial learning rate

e.g.  $\alpha_0 = 0.2$ ,  $\text{decay\_rate} = 1$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04

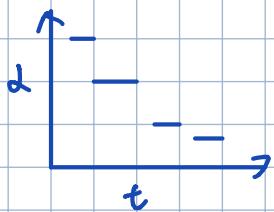
## Other learning rate decay method

$$\alpha = 0.95^{\text{epoch\_num}} \cdot \alpha_0$$

exponentially decay

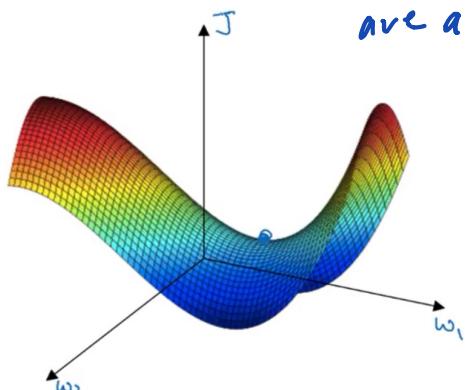
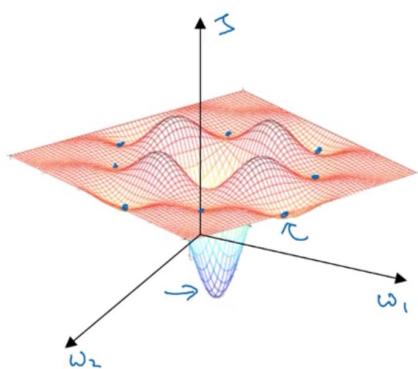
$$\alpha = \frac{K}{\sqrt{\text{epoch\_num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{K}{\sqrt{t}} \alpha_0$$

or discrete decay



## The problem of local optima

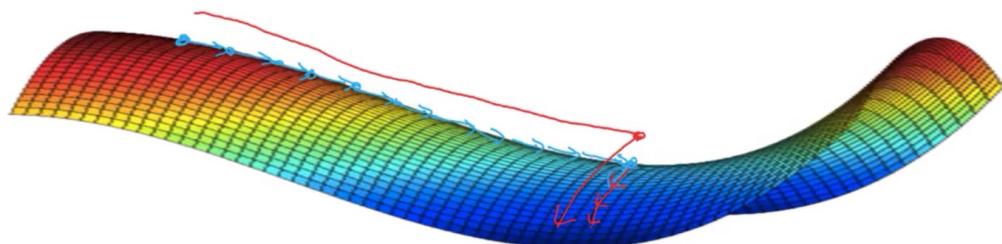
### Local optima in neural networks



many "local optima" in NN  
are actually saddle points

i.e. if 20,000 features  
then for it to be  
local minima,  
all 20,000 directions  
has to be  $\cup$

### Problem of plateaus → can really



## Week 3:

### Hyperparameter tuning, Batch Normalisation and Programming framework

#### Tuning Process

Some guidelines to systematically organise hyperparameter tuning process.

- 1.  $\lambda$   $\rightarrow$  import
- 2.  $\beta_1 \sim 0.9$   $\beta_2 \sim 0.999$   $\epsilon \sim 10^{-8}$
- 3.  $\beta_1, \beta_2, \epsilon$  (Adam)  $\rightarrow$  2nd import
- 4. learning rate decay  $\rightarrow$  3rd import
- 5. layers
- 6. hidden units for different layers
- 7. learning rate decay
- 8. mini batch size

#### Grid Search

hyperparameter 2

hyperparameter 1	- - - - .
	- - - - -
	- - - - ^
	- - - -

works ok when a hyperparameter was relatively small

#### Random Search (for deep learning)

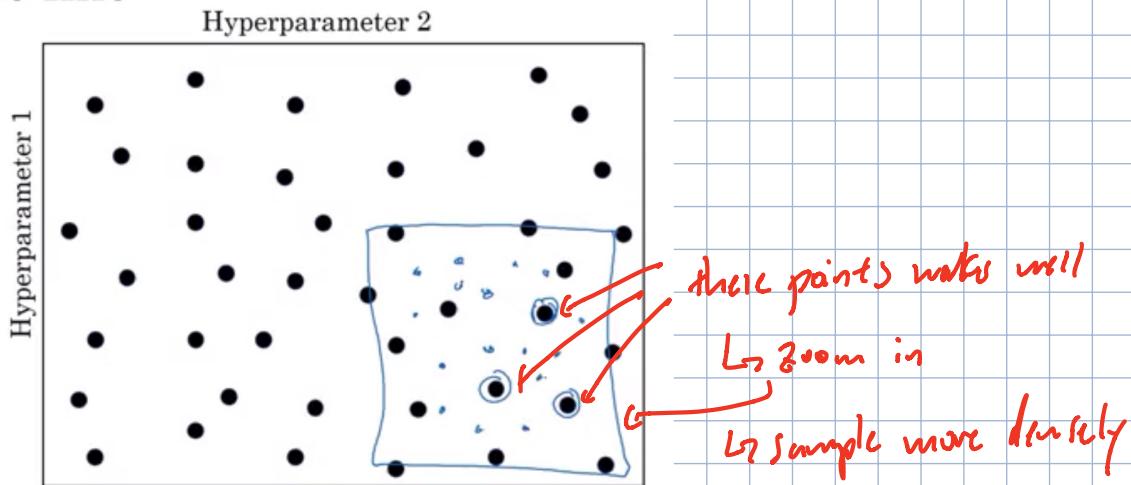
hyperparameter 2

hyperparameter 1	.
	.
	..
	- -

works better because not all hyperparameters has the same importance.

Another common practice is to use coarse-to-fine sampling scheme

## Coarse to fine



## Using an appropriate scale to pick hyperparameters

Sampling hyperparameters at random  $\neq$  Sampling uniformly at random over the range of valid values.  $\rightarrow$  it is important to pick the right scale.

e.g. picking hyperparameters at random

say we try to choose # hidden units,  $N^{(l)}$ , for layer  $l$

we think a good range of value is  $50 \leftrightarrow 100$ :

$$N^{(l)} = 50, \dots, 100$$

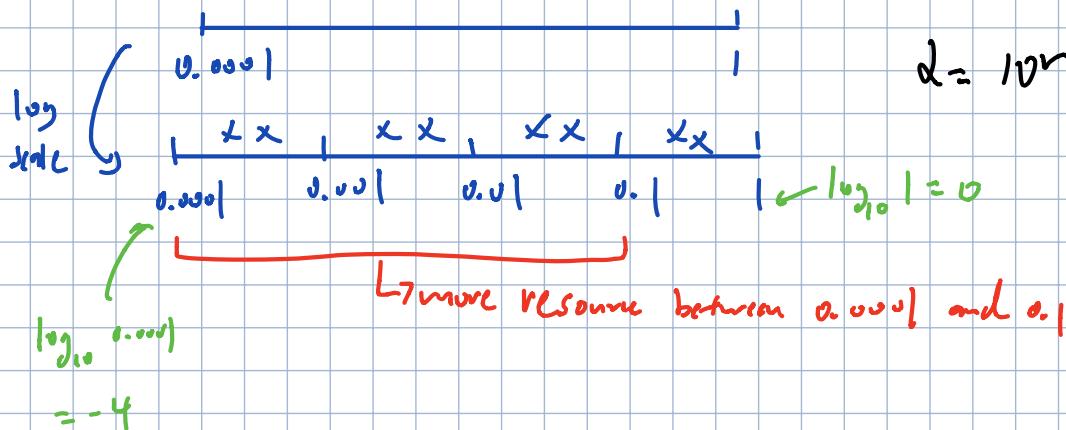
$$\underbrace{ \quad x \quad x \quad x \quad x \quad x \quad x \quad }_{50} \quad \underbrace{\quad x \quad x \quad x \quad}_{100}$$

over layers  $L: 2, 3, 4, 5$

} examples that uniform random sampling are reasonable but this is not true for all hyperparameters

Say we searching for learning rate,  $\alpha$

$$\alpha = 0.0001, \dots, 1$$



$$r = -4 \times \text{np.random.rand}()$$

$$\alpha = 10^r$$

Treky case: Sampling the hyperparameter ( $\beta$ ) for exponentially weighted averages.

$$\beta = 0.9, \dots, 0.999$$

↓                          ↓  
average                average  
last 10 values      last 1,000 values

$$\text{Can think of } 1-\beta = 0.1, \dots, 0.001$$

$$\begin{array}{c} \hline 0.1 & 0.01 & 0.001 \\ \hline 10^{-1} & & 10^{-3} \end{array} \rightarrow r \in [-3, -1]$$

Reason:

$$\beta: 0.9 \rightarrow 0.999$$

$$\hookrightarrow 1-\beta = 10^r \rightarrow \beta = 1-10^r$$

hardly any change in result

$$0.999 \rightarrow 0.9995$$



## Hyperparameters tuning in practice: Pandas vs Caviar

### ① Build 1 model

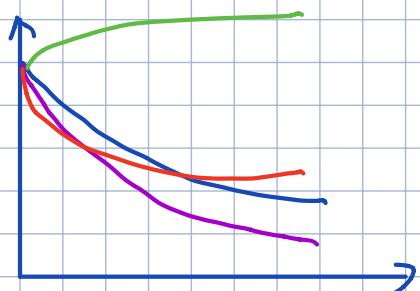
large dataset, but limited computational resources



watching performance and partially nudging  
learning rate up or down

"Panda"

### ② Train many model in parallel



"Caviar"

## Batch Normalisation

### Normalising activations in a network

Makes hyperparameter search problem much easier, and makes  
NN much more robust

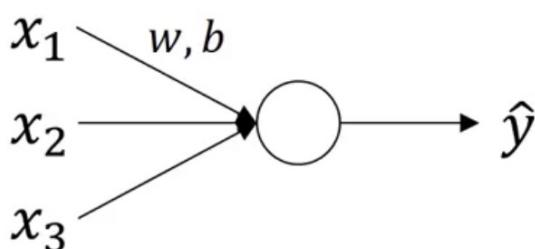
when training a model (e.g. logistic  
regression), normalising the input feature  
can speed up learnings

$$\text{recall: } M = \frac{1}{m} \sum_i X^{(i)}$$

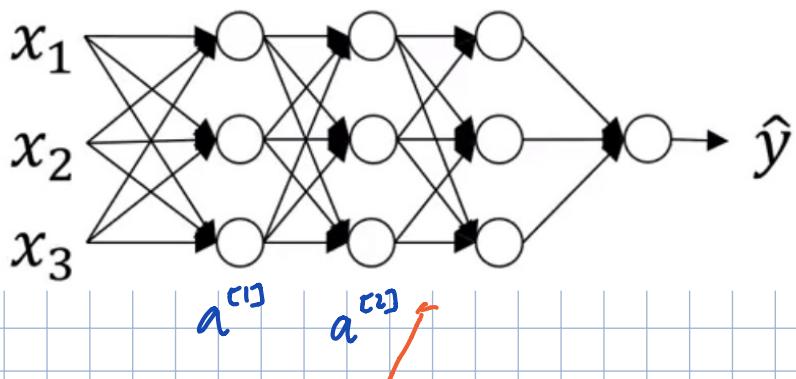
$$X = X - M \quad \text{demographic}$$

$$\sigma^2 = \frac{1}{m} \sum_i X^{(i)}_i$$

$$X = X / \sigma^2$$



How about a deeper model?



if we want train parameters

$$w^{(1)}, b^{(1)}$$

↳ normalise  $a^{(2)}$  to make the training of  $w^{(2)}, b^{(2)}$  more efficient

For any hidden layer, can we normalise the value of  $a^{(l)}$  so as to train  $w^{(l+1)}, b^{(l+1)}$  fast?

↳ Batch Normalisation!

↳ we will be normalising  $z^{(l)}$  (instead of  $a^{(l)}$ )

### Implement Batch Norm

Given some intermediate values in NN:  $z^{(1)}, \dots, z^{(m)}$  for some layer l  
 $z^{(l)(i)}$

Compute  $M = \frac{1}{m} \sum_i z^{(i)}$

(for some layer l, omitting (l) notation)

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

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \text{numerical stability}$$

every component of  $z$   
mean = 0  
var = 1

we don't want the hidden unit always have mean 0 and var=1

↳ have different distribution

not only has input features,  
but also activations  $a^{(0)} \dots$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

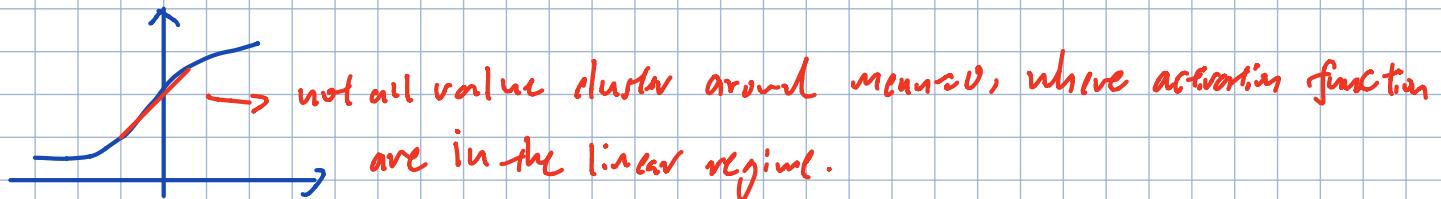
learnable parameters from model

from e.g. gradient descent  
we will update  $\gamma$  and  $\beta$   
just as  $w$  and  $b$

$\gamma$  and  $\beta$  allow us to set the mean of  $\tilde{z}$  to any value:

if  $\gamma = \sqrt{\sigma^2 + \epsilon}$  } invert back to  $z^{(i)}$  (from  $z_{\text{norm}}^{(i)}$ )  
 $\beta = \mu$  } i.e.  $\tilde{z}^{(i)} = z^{(i)}$

we will use  $\tilde{z}^{(l+1)}$ , instead of  $z^{(l+1)}$ , for the latter we write



### Fitting Batch Norm into a NN

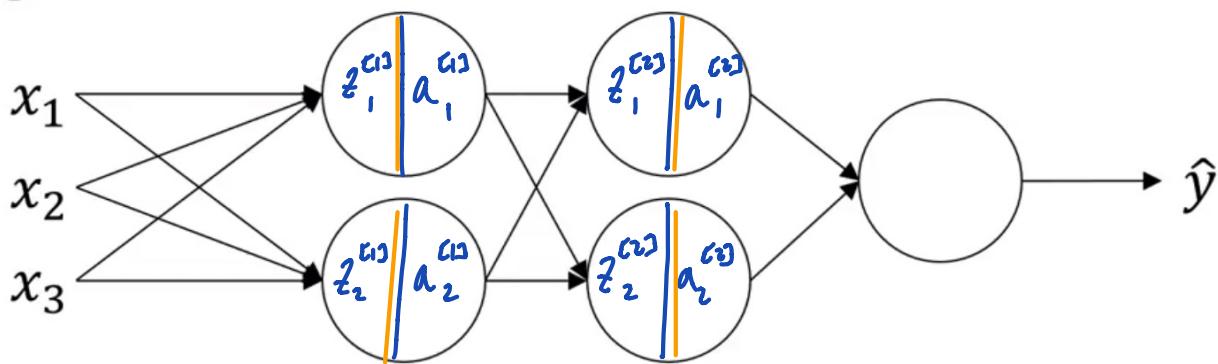


Diagram illustrating the flow of data through a neural network layer, showing the application of Batch Normalization (BN) at each stage.

The input  $x$  is processed by weights  $w^{(1)}$  and bias  $b^{(1)}$  to produce the pre-activations  $z^{(1)}$ .

The pre-activations  $z^{(1)}$  are then processed by Batch Normalization (BN) using scale factor  $\gamma^{(1)}$  and shift  $\beta^{(1)}$  to produce the normalized features  $\tilde{z}^{(1)}$ .

The normalized features  $\tilde{z}^{(1)}$  are then processed by weights  $w^{(2)}$  and bias  $b^{(2)}$  to produce the pre-activations  $z^{(2)}$ .

The pre-activations  $z^{(2)}$  are then processed by BN using scale factor  $\gamma^{(2)}$  and shift  $\beta^{(2)}$  to produce the final output  $\hat{y}$ .

Parameters of our network:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]}$   
 $\beta^{[1]}, \gamma^{[2]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$

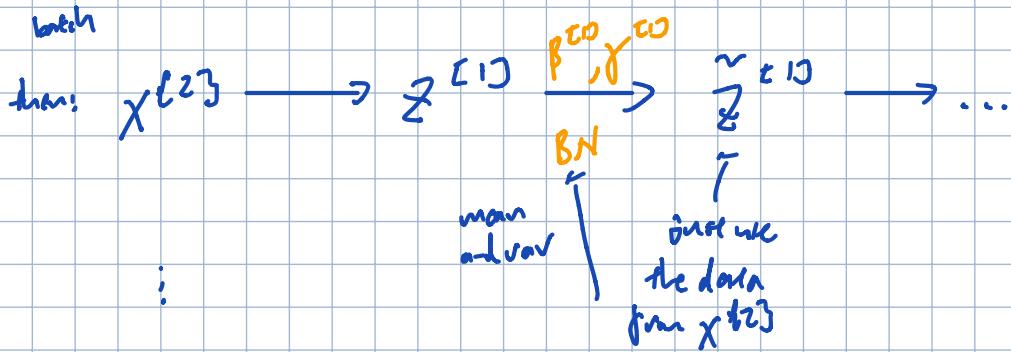
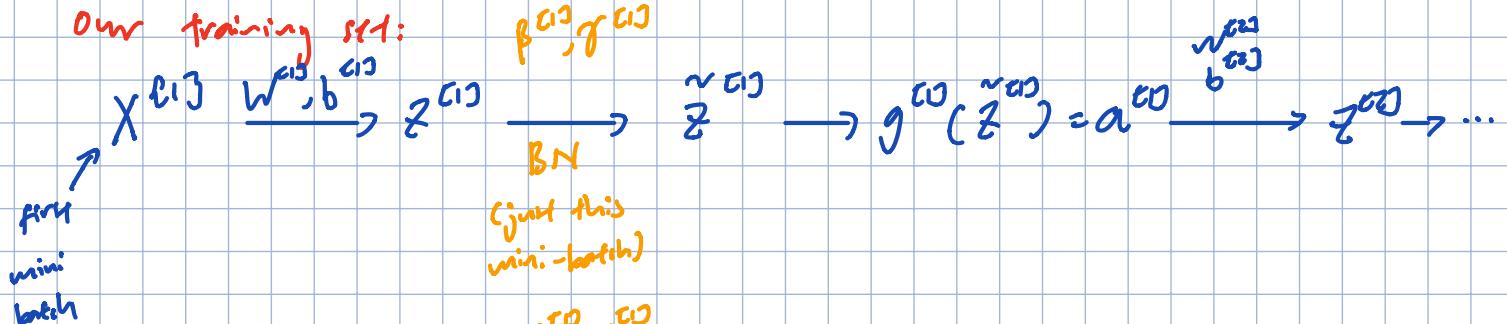
$\hat{\beta}$   
 nothing to do with momentum hyperparameter  $\beta$   
 RMSprop  
 Adam

We can still use gradient descent, but need

$$\lambda\beta^{[L]}, \quad \beta^{[L]} = \beta^{[L]} - \alpha\lambda\beta^{[L]}$$

In practice, batch norm is usually applied with mini-batch of

Our training set:  $\beta^{[1]}, \gamma^{[1]}$



Parameters  $W^{[L]}, b^{[L]}, \beta^{[L]}, \gamma^{[L]}$

$$\hookrightarrow Z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]} \rightarrow \text{gets canceled out during Batch Norm}$$

(  
↳ home

$$Z^{[L]} = W^{[L]} a^{[L-1]}$$

(  
↳  $Z^{[L]}$ )

$$( \quad \tilde{Z}^{[L]} = \gamma^{[L]} Z_{\text{norm}}^{[L]} + \beta^{[L]}$$

↳ any constant added here doesn't change anything during Batch Norm  
 (cancel out during mean subtraction)

for 1 example

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

$$b^{[t]} : (n^{[t]}, 1)$$

$$\beta^{[t]} : (n^{[t]}, 1)$$

$$\gamma^{[t]} : (n^{[t]}, 1)$$

### Implementing Gradient Descent with Batch Norm (mini-batch)

for  $t=1 \dots \text{num\_mini-batch}$

compute forward prop on  $X^{[t]}$

In each hidden layer, use BN to replace  $z^{[t]}$  with  $\tilde{z}^{[t]}$

use back prop to compute  $dW^{[t]} - \cancel{d\tilde{z}^{[t]}}$ ,  $d\beta^{[t]} - d\gamma^{[t]}$

$$\text{update } W^{[t]} := W^{[t]} - \alpha dW^{[t]}$$

$$\beta^{[t]} := \beta^{[t]} - \alpha d\beta^{[t]} \quad \left. \right\} \text{gradient descent}$$

$$\gamma^{[t]} := \gamma^{[t]} - \alpha d\gamma^{[t]}$$

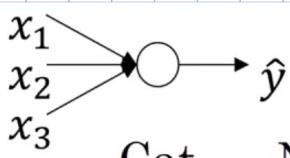
also works with GD with momentum, RMSprop, Adam

### Why does Batch Norm work?

First reason: similar to normalize input in logistic regression for GD.

Just that this is applied to further hidden nodes

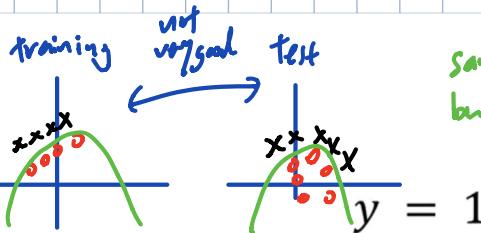
Second reason: makes weights flatter in NN more robust to changes to weights in earlier layers of NN



Cat  
 $y = 1$



Non-Cat  
 $y = 0$



same func might work  
but algorithm won't learn this  
by just looking at  
train data



Covariate shift

$X \rightarrow Y$

$T$   
if change  
 $\Rightarrow$  weight need to  
retrain

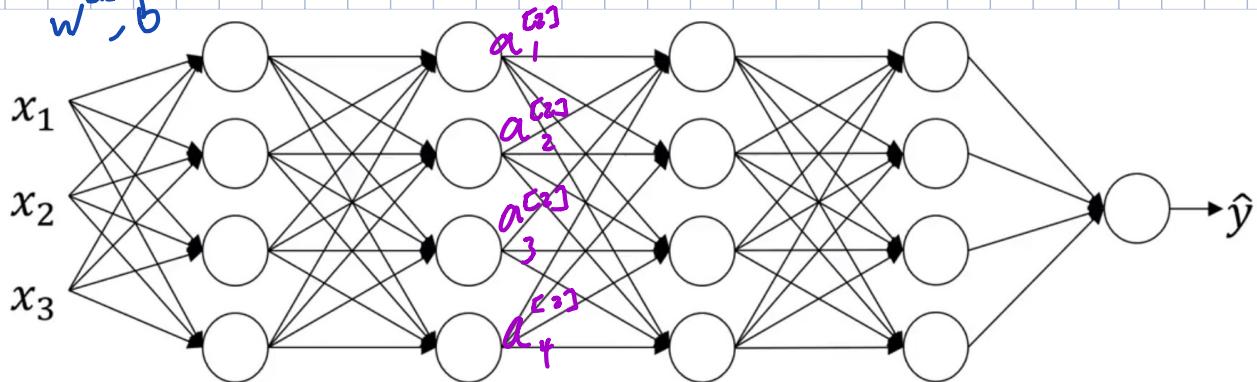
↑  
train on  
black cats

don't do  
very well

↑  
predict on  
color cats

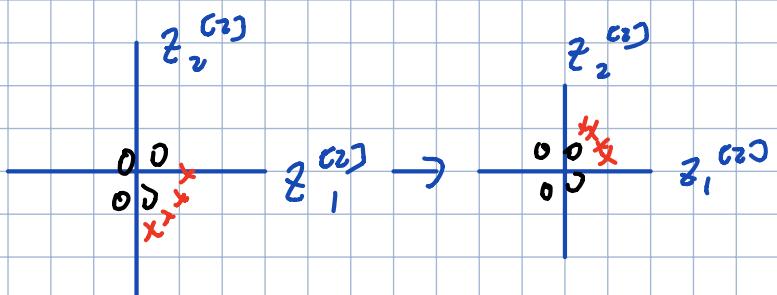
This is a problem for NN : let's look at 3rd layer below

$$w^{(0)}, b^{(0)} \quad w^{(1)}, b^{(1)} \quad w^{(2)}, b^{(2)} \quad w^{(3)}, b^{(3)} \downarrow \quad w^{(4)}, b^{(4)}$$



3rd layer job: take  $a^{(2)}$   $\xrightarrow{\text{map}}$   $\hat{y}$ . But  $a^{(2)}$  are from layer 2, which changes  
from 1st layer  $\Rightarrow$  i.e.  $a^{(2)}$  changes all the time  $\rightarrow$  covariate shift

Batch norm : reduces the amount that the distribution of these hidden unit values shifts around.



even if exact value of  
 $z_1^{[l]}, z_2^{[l]}$  changes, its  
 mean/var will at least  
 stay the same

$\swarrow$   
 allows input to the next more stable

OR

It weakens the coupling between early layers parameters how to do with later layers, so it allows each layer to learn by itself.

third reason: slight regularization effect

## Batch Norm as regularization

$$\rightarrow \{x_i\} : z^{[l]}$$

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values  $z^{[l]}$  within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

$$z^{[l]} \rightarrow \tilde{z}^{[l]}$$

increase minibatch size

reduce noise

reduce regularization budget

## Batch Norm at test time

At test time, may not have minibatch, maybe one at a time. we have to adapt our NN to do that.

Local Batch Norm has this issue, as it changes forward prop (i.e. prediction)

# Batch Norm at test time

*(during training)*

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

*(computed using a single mini-batch)*

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

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

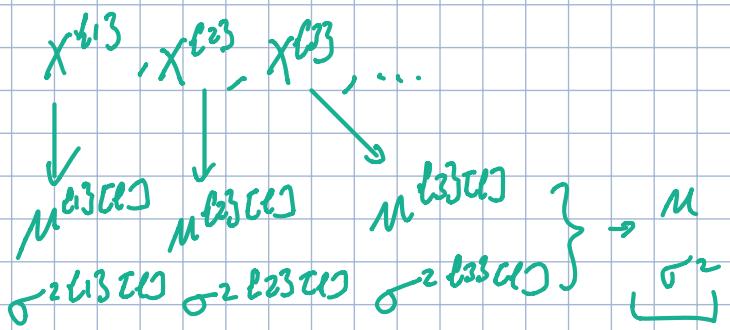
$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

In test time, only 1 example of a time to predict.

In test time, come up with some separate estimate of  $\mu$  and  $\sigma^2$

↳ estimate using exponentially weighted average (across mini batch)

e.g. for layer  $l$



than when making a prediction!

$$z_{\text{norm}} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z} = \gamma z_{\text{norm}} + \beta \Rightarrow \hat{y}$$

*from NN*

OR, you can run the entire training set, then get average/final  $M, \sigma^2$

## Multi-class classification

### Softmax Regression

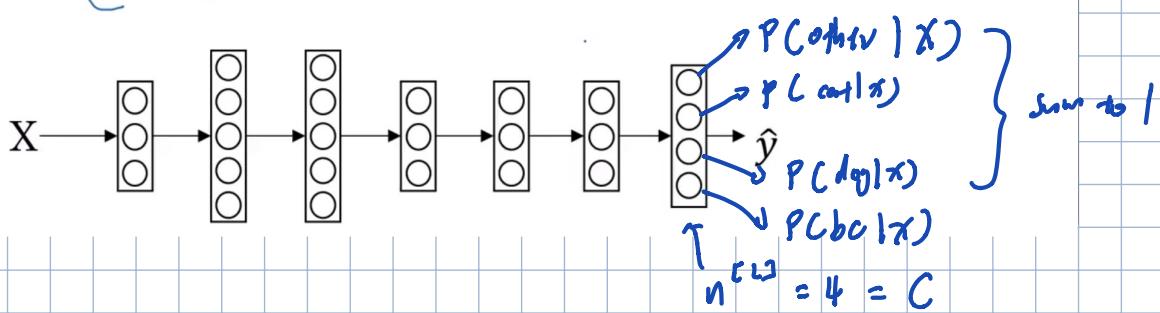
It is a generalised form of logistic regression

Recognizing cats, dogs, and baby chicks, other



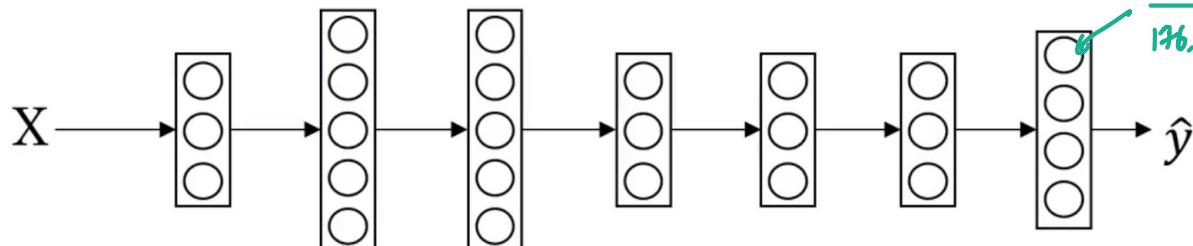
3      1      2      0      3      2      0      1

$$C = \# \text{classes} = 4 \quad (0, \dots, 3)$$



$\hat{y}$  is  $(4, 1)$  for 1 example

### Softmax layer



At the final layer L:

$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

Activation function: temporary variable  $t = e^{(z^{[L]})}$  element-wise share for 1 example  $z^{[L]} = (4, 1)$ , so it's also  $(4, 1)$

$$\rightarrow t = e^{(z^{[L]})}$$

$$\rightarrow a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{i=1}^4 t_i} \quad \text{also } (4, 1)$$

earliest element

$$\hookrightarrow a_i^{[L]} = \frac{t_i}{\sum_{i=1}^4 t_i}$$

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \Rightarrow t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}$$

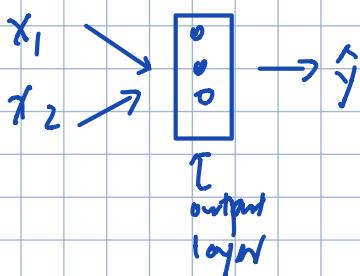
$$\hookrightarrow \sum_{i=1}^4 t_i = 176.3 \Rightarrow a^{[L]} = t / 176.3 = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

softmax activation function:

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

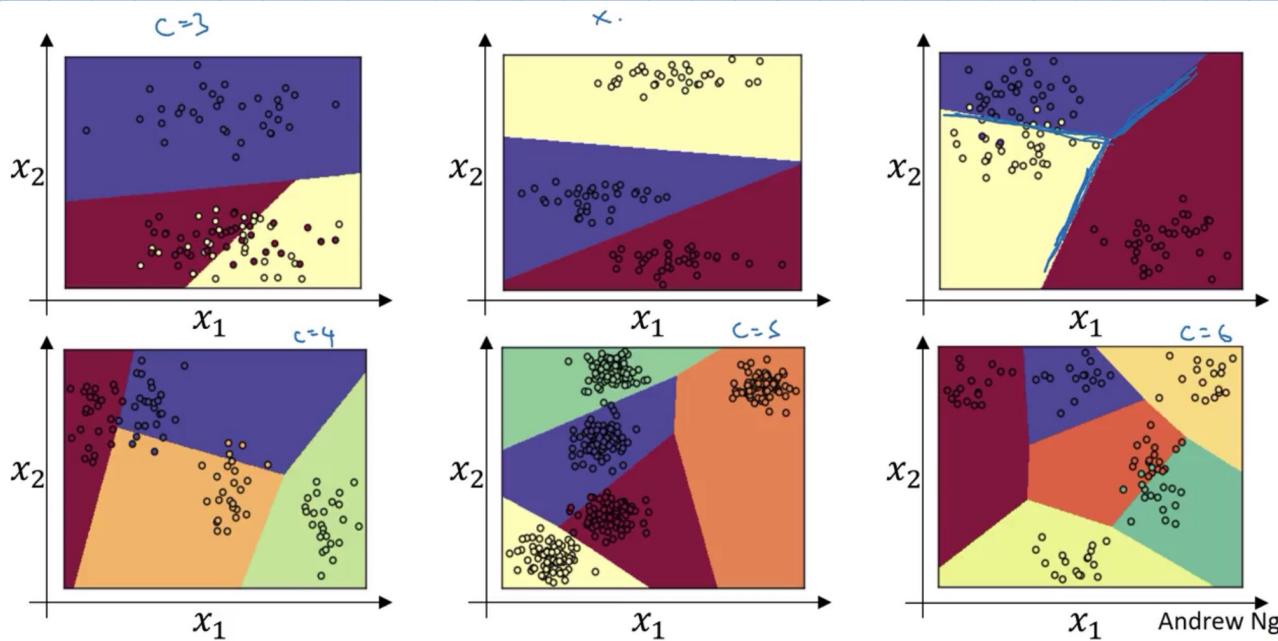
$$g^{[L]}_{(4,1)} = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

A NN with no hidden layers!



$$z^{[L]} = W^{[L]} x + b^{[L]}$$

$$a^{[L]} = \hat{y} = g(z^{[L]})$$



## Training a softmax classifier

$$(4,1) \quad z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$$C=4$$

$$\begin{aligned} a^{[L]} &= \begin{bmatrix} e^5/(e^5 + e^2 + e^{-1} + e^3) \\ e^2/(e^5 + e^2 + e^{-1} + e^3) \\ e^{-1}/(e^5 + e^2 + e^{-1} + e^3) \\ e^3/(e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \\ g^{[L]}(z^{[L]}) &= \text{"softmax"} \quad \text{"hard max"} \\ &= \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

softmax activation function generalizes the logistic activation function to C classes (rather than 2 classes) (so it is still a linear classifier)

If C=2, softmax reduces to logistic regression

$$\text{i.e. } a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix} \xleftarrow{\text{redundant}} \text{1 value} \Rightarrow \text{logistic regression}$$

How we train a NN with a softmax output layer?

$$\text{e.g. true label } y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ cat} \quad \text{prediction } \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} = a^{[L]} \quad C=4$$

$$\text{Loss function } L(\hat{y}, y) = - \sum_{j=1}^4 y_j \log \hat{y}_j \quad (\text{1 example})$$

$$\text{here } y_1 = y_3 = y_4 = 0 \\ y_2 = 1 \quad \longrightarrow \quad -y_2 \log \hat{y}_2 = -\log \hat{y}_2$$

i.e. it algorithm try to make  $L(\hat{y}, y)$  small (through gradient descent)

↳ make  $-\log \hat{y}_2$  small

↳ make  $\hat{y}_2$  big

this loss function look at true label, then try to make the corresponding prediction probability big.

cost function  $J$  on the entire training set

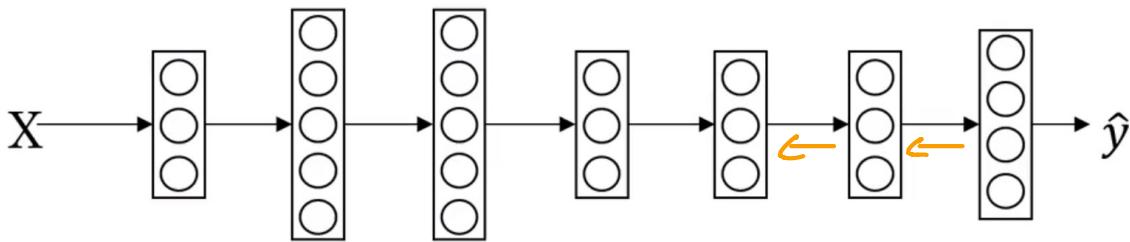
$$J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

vectorized:

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \dots$$

$$\hat{Y} = [\hat{y}^{(1)} \ \hat{y}^{(2)} \ \dots \ \hat{y}^{(m)}] = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.6 \end{bmatrix} \dots$$

## Gradient descent with softmax



Back prop (with softmax at layer  $L$ ):

$$d\hat{z}^{[L]} = \hat{y} - y, \text{ then starts the back prop}$$

$$\frac{\partial J}{\partial z^{[L]}}$$

$$z^{[L]} \rightarrow a^{[L]} = \hat{y}$$

(4,1)

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

## Deep Learning frameworks

# Deep learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

### Choosing deep learning frameworks

- Ease of programming (development and deployment)
- Running speed
- Truly open (open source with good governance)

## TensorFlow

### Motivation

e.g. minimize  $J(w) = w^2 - 10w + 25$

$$(w-5)^2 \rightarrow w=5 \text{ is minimum } J$$

```
In [1]: import numpy as np
import tensorflow as tf
```

```
In [5]: w = tf.Variable(0,dtype=tf.float32)
cost = tf.add(tf.add(w**2,tf.multiply(-10.,w)),25)
cost = w**2 - 10*w + 25
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
init = tf.global_variables_initializer()
session = tf.Session()
session.run(init)
print(session.run(w))
```

0.0

```
In [6]: session.run(train)
print(session.run(w))
```

0.1

```
In [7]: for i in range(1000):
    session.run(train)
print(session.run(w))
```

4.9999

w has to be declared as  
tf.variable

Starts a session

initialise global variable

→ run 1 step of gradient descent

# Code example

```
import numpy as np
import tensorflow as tf

coefficients = np.array([[1], [-20], [25]])

w = tf.Variable([0], dtype=tf.float32)
x = tf.placeholder(tf.float32, [3,1])
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]    # (w-5)**2
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()

session = tf.Session()                         with tf.Session() as session:
session.run(init)                            session.run(init)
print(session.run(w))                        print(session.run(w))

for i in range(1000):
    session.run(train, feed_dict={x:coefficients})
print(session.run(w))
```