**Course 2 Improving Deep Neural Networks:**

**Hyperparameter tuning, Regularization and Optimization**

Study Notes by Jane Huang

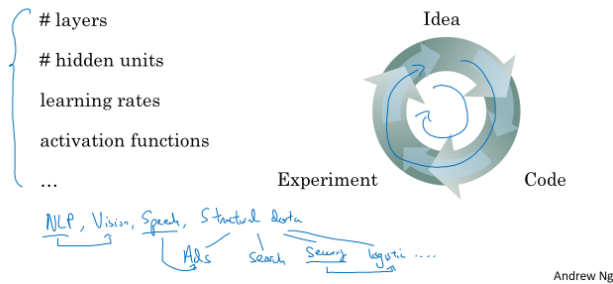08172018

# Table of Contents

**Week 1** **Practical aspects of Deep Learning**

**1.1 Train/dev/test sets**

Applied ML is a highly iterative process

# layers                    Idea
# hidden units
learning rates
activation functions
...              Experiment        Code

NLP, Vision, Speech, Structural data
Ads    search   Security   logistic ....

Andrew Ng

Train/dev/test sets

Data        dev    test
        training set      - Hold-out cross validation    test
                          - Development set "dev"

Prev era:  70/30 %      60/20/20 %
           100 - 1000 - 10000

Big data:   1,000,000      10,000      10,000
            98 / 1 / 1 %.
            99.5 { 75 / 75
                   .4 / -1 %.

**1.2 Bias and Variance**



high bias          "just right"          high variance
underfitting                             overfitting

**1.3 Regularization**

If you use L1 regulariztion, w will be sparse. If you use L2 regularization, "weight decay"

2

$$\min_{w,b} J(w,b)$$

$w \in \mathbb{R}^{n_x}, \quad b \in \mathbb{R} \qquad \lambda = \text{regularization parameter}$

lambda        lambd

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

omit

$L_2$ regularization $\qquad \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$

$L_1$ regularization $\qquad \dfrac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \dfrac{\lambda}{2m} \|w\|_1$

$w$ will be sparse

# Neural network

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

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

"Frobenius norm" $\qquad \|\cdot\|_2^2 \qquad \|\cdot\|_F^2$

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

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

"Weight decay" $\qquad w^{[l]} := w^{[l]} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right]$

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

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

Andrew N

"Frobenius norm"

# How does regularization prevent overfitting?



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

$w^{[l]} \approx 0$

high bias          "just right"          high variance

# How does regularization prevent overfitting?



Every layer ~ Linear

If the regularization becomes very large, the parameters W very small, so Z will be relatively small, kind of ignoring the effects of b for now, so Z will be relative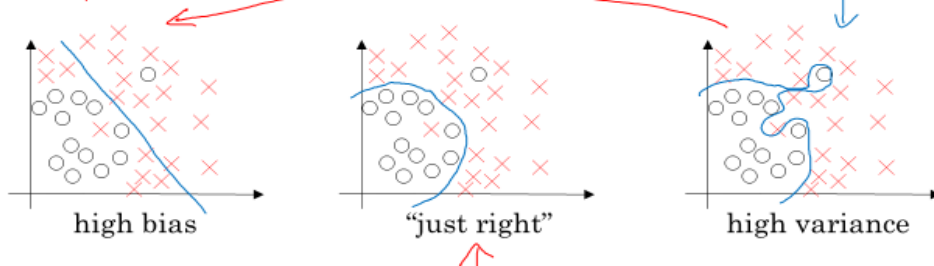ly small or, really, I should say it takes on a small range of values. And so the activation function if is tanh, say, will be relatively linear. And so your whole neural network will be computing something not too far from a big linear function which is therefore pretty simple function rather than a very complex highly non-linear function. And so is also much less able to overfit.

## 1.4 Dropout regularization



After dropout, you end up with a much smaller, much diminished network.

D3 is a Boolean array

In this example, dropout=0.2, keep-prob=1-dropout=0.8

# Implementing dropout ("Inverted dropout")

Illustrate with layer $l = 3$.   keep-prob = 0.8 ————   0.2

$\rightarrow$  d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keep-prob

a3 = np.multiply(a3, d3)        # a3 *= d3.

$\rightarrow$  a3 /= ~~0.8~~ keep-prob  $\leftarrow$

50 units.  $\leadsto$   10 units shut off

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

J          $\uparrow$ reduced by 20%.       Test

/= 0.8

## Making predictions at test time

$a^{[0]} = X$

No drop out.

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

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

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

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

/= keep-prob

$\downarrow$
$\hat{y}$

Don't implement dropout at test time, otherwise it is adding noise

No need to add extra scaling parameter

## Why does drop-out work?

Intuition: <u>Can't rely on any one feature, so have to spread out weights.</u> $\leadsto$ Shrink weights. $L_2$

$\rightarrow$ Computer Vision

Domsik: J

$x_1$
$x_2$
$x_3$

keep-prob = 1   0.7   0.5 $\leftarrow$

0.9   $W^{[3]}$

#iteration

Drop out has similar effect as L2 regularization

But for layers where you're more worried about over-fitting, really the layers with a lot of parameters, you can set the key prop to be smaller to apply a more powerful form of drop out.

5

It's kind of like cranking up the regularization parameter lambda of L2 regularization where you try to regularize some layers more than others. Normally no dropout at input layer.

**1.5 Other regularization methods**

Data augmentation



Early stopping

Orthogonalization

- Optimizing the cost function
    - Gradient
- Not overfit
    - Regularization

And the advantage of early stopping is that running the gradient descent process just once, you get to try out values of small w, mid-size w, and large w, without needing to try a lot of values of the L2 regularization hyperparameter lambda. But I (Andrew) personally prefer to just use L2 regularization and try different values of lambda. That's assuming you can afford the computation to do so. But early stopping does let you get a similar effect without needing to explicitly try lots of different values of lambda. So you've now seen how to use data augmentation as well as if you wish early stopping in order to reduce variance or prevent over fitting your neural network.

## 1.6 Normalizing inputs

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



Substract mean:

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

$x := x - \mu$

Normalize variance:

$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} ** 2$   ← element-wise

$x /= \sigma^2$

Use same $\mu$ $\sigma^2$ to normalize test set.

# Why normalize inputs?

$w_1$   $x_1 : 1 \dots 1000$ ←
$w_2$   $x_2 : 0 \dots 1$ ←

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

$-1 \dots 1$

Unnormalized:



Normalized:



$x_1 : 0 \dots 1$
$x_2 : -1 \dots 1$
$x_3 : 1 \dots 2$

If you use unnormalized input features, it's more likely that your cost function will look like this, it's a very squished out bowl, very elongated cost function, where the minimum you're trying to find is maybe over there. But if your features are on very different scales, say the feature X1 ranges from 1 to 1,000, and the feature X2 ranges from 0 to 1, then it turns out that the ratio or the range of values for the parameters w1 and w2 will end up taking on very different values. And so maybe these axes should be w1 and w2, but I'll plot w and b, then your cost function can be a very elongated bowl like that. So if you part the contours of this function, you can have a very elongated function like that. Whereas if you normalize the features, then your cost function will on average look more symmetric. And if you're

7

running gradient descent on the cost function like the one on the left, then you might have to use a very small learning rate because if you're here that gradient descent might need a lot of steps to oscillate back and forth before it finally finds its way to the minimum. Whereas if you have a more spherical contours, then wherever you start gradient descent can pretty much go straight to the minimum. You can take much larger steps with gradient descent rather than needing to oscillate around like like the picture on the left. Of course in practice w is a high-dimensional vector, and so trying to plot this in 2D doesn't convey all the intuitions correctly. But the rough intuition that your cost function will be more round and easier to optimize when your features are all on similar scales.

**1.7 Vanishing / Exploding gradient**



If linear activation function, for example, g(z)=z

If w_l>1 exploding gradients

If w_l<1 vanishing gradients

So the intuition I hope you can take away from this is that at the weights W, if they're all just a little bit bigger than one or just a little bit bigger than the identity matrix, then with a very deep network the activations can explode. And if W is just a little bit less than identity. So this maybe here's 0.9, 0.9, then you have a very deep network, the activations will decrease exponentially. And even though I went through this argument in terms of activations increasing or decreasing exponentially as a function of L, a similar argument can be used to show that the derivatives or the gradients the computer is going to send will also increase exponentially or decrease exponentially as a function of the number of layers. And this makes training difficult, especially if your gradients are exponentially smaller than L, then gradient descent will take tiny little steps. It will take a long time for gradient descent to learn anything.

To summarize, you've seen how deep networks suffer from the problems of vanishing or exploding gradients. In fact, for a long time this problem was a huge barrier to training deep neural networks. It turns out there's a partial solution that doesn't completely solve this problem but it helps a lot which is careful choice of how you initialize the weights.

### 1.8 Weight Initialization for Deep Networks

A partial solution to vanishing/exploding gradient, doesn't solve it entirely but helps a lot, is better or more careful choice of the random initialization for your neural network.



Xavier Initialization

### 1.9 Numerical approximation of gradients

# Checking your derivative computation

$f(\theta) = \theta^3$



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

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

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

approx error: $0.0001$
(prev slide: $3.0301$. error: $0.03$)

$$\left( f'(\theta) = \lim_{\epsilon \to 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \right) \quad O(\epsilon^2) \quad \left| \quad \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \quad \text{error: } O(\epsilon) \right.$$

$0.01$     $0.01$
$0.0001$

Andrew Ng

**1.10 Gradient Checking**

# Gradient check for a neural network

Take $W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]}$ and reshape into a big vector $\theta$.

concatenate

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

Take $dW^{[1]}, db^{[1]}, \ldots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

concatenate

Is $d\theta$ the gradient of $J(\theta)$?

- You gotta take all of these Ws and reshape them into vectors, and then concatenate all of these things, so that you have a giant vector theta. The cost function is a function of theta, therefore a function of the Ws and Bs.
- You gotta take all of these dWs and reshape them into vectors, and then concatenate all of these things, so that you have a giant vector dtheta.

# Gradient checking (Grad check)

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

for each $i$:

$$\rightarrow d\theta_{appx}[i] = \frac{J(\theta_1, \theta_2 \dots, \theta_i + \varepsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \varepsilon, \dots)}{2\varepsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \qquad d\theta_{appx} \overset{?}{\approx} d\theta$$

$i$

Check $\quad \dfrac{\| d\theta_{appx} - d\theta \|_2}{\| d\theta_{appx} \|_2 + \| d\theta \|_2}$

$$\varepsilon = 10^{-7}$$

$x \quad \boxed{\dfrac{10^{-7}}{10^{-5}} \quad - \; great! \; \leftarrow}$

$\rightarrow 10^{-3} \; - worry. \; \leftarrow$

# Gradient checking implementation notes

- Don't use in training – only to debug

$$d\theta_{appx}[i] \longleftrightarrow d\theta[i]$$
$\qquad\qquad \uparrow \;\uparrow \qquad\qquad \uparrow$

- If algorithm fails grad check, look at components to try to identify bug

$$db^{[l]}_{k} \qquad dw^{[l]}_{k}$$

- Remember regularization.

$$J(\theta) = \frac{1}{m} \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_l \| w^{[l]} \|_F^2$$

$d\theta = $ grad of $J$ w.r.t. $\theta$

- Doesn't work with dropout. $\quad J \qquad keep\text{-}prob = 1.0$

- Run at random initialization; perhaps again after some training.

$$W, b \approx 0$$

You can do the grad check by setting keep_prob=1.0. If it is correct, then turn on drop out.

**Week 2** Optimization algorithms

**2.1 Mini-batch gradient descent**

## Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on $m$ examples.

$$X = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} | x^{(1001)} \ \dots \ x^{(2000)} | \ \dots \ | \ \dots \ x^{(m)}]$$

$(n_x, m)$

$X^{\{1\}} \ (n_x, 1000) \qquad X^{\{2\}} \ (n_x, 1000) \ \dots \qquad X^{\{5000\}} \ (n_x, 1000)$

$$Y = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} | y^{(1001)} \ \dots \ y^{(2000)} | \ \dots \ | \ \dots \ y^{(m)}]$$

$(1, m)$

$Y^{\{1\}} \ (1, 1000) \qquad Y^{\{2\}} \ (1, 1000) \qquad Y^{\{5000\}} \ (1, 1000)$

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

$5,000$ mini-batches of $1,000$ each

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

$x^{(i)}$

$z^{[\ell]}$

$X^{\{t\}}, Y^{\{t\}}$

Andrew Ng

With the implementation of gradient descent on your whole training set, what you have to do is, you have to process your entire training set before you take one little step of gradient descent. And then you have to process your entire training sets of five million training samples again before you take another little step of gradient descent. So, it turns out that you can get a faster algorithm if you let gradient descent start to make some progress even before you finish processing your entire, your giant training sets of 5 million examples.

Let's say that you split up your training set into smaller, little baby training sets and these baby training sets are called mini-batches.

Use curly bracket to index mini-batch

# Mini-batch gradient descent

repeat {

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

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

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$\vdots$$
$$A^{[L]} = g^{[L]}(Z^{[L]})$$

    } Vectorized implementation (1000 examples)

    Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{\ell} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{\ell} \|w^{[\ell]}\|_F^2$.

    Backprop to compute gradients w.r.t $J^{\{t\}}$ (using $(X^{\{t\}}, Y^{\{t\}})$)

$$W^{[\ell]} := W^{[\ell]} - \alpha \, dW^{[\ell]}, \quad b^{[\ell]} := b^{[\ell]} - \alpha \, db^{[\ell]}$$

}
}

*1 step of gradient descent using $X^{\{t\}}, Y^{\{t\}}$ (of 1m ≤ 1000)*

$X, Y$

*for $X^{(i)}, y^{(i)}$.*

*"1 epoch"* — pass through training set.

Andrew Ng

The code I have written down here is also called doing one epoch of training and epoch is a word that means a single pass through the training set. Whereas with batch gradient descent, a single pass through the training allows you to take only one gradient descent step. With mini-batch gradient descent, a single pass through the training set, that is one epoch, allows you to take 5,000 gradient descent steps. Now of course you want to take multiple passes through the training set which you usually want to, you might want another for loop for another while loop out there. So you keep taking passes through the training set until hopefully you converge with approximately converge. When you have a lot training set, mini-batch gradient descent runs much faster than batch gradient descent and that's pretty much what everyone in Deep Learning will use when you're training on a large data set.

## Training with mini batch gradient descent

**Batch gradient descent**



**Mini-batch gradient descent**



$X^{\{1\}}, Y^{\{1\}}$
$X^{\{2\}}, Y^{\{2\}}$
$\vdots$
$J^{\{t\}}$

Plot $J^{\{t\}}$ computed using $X^{\{t\}}, Y^{\{t\}}$

Andrew Ng

- With batch gradient descent on every iteration you go through the entire training set and you'd expect the cost to go down on every single iteration. So if we've had the cost function j as a

function of different iterations it should decrease on every single iteration. And if it ever goes up even on iteration then something is wrong. Maybe you're running ways to big.

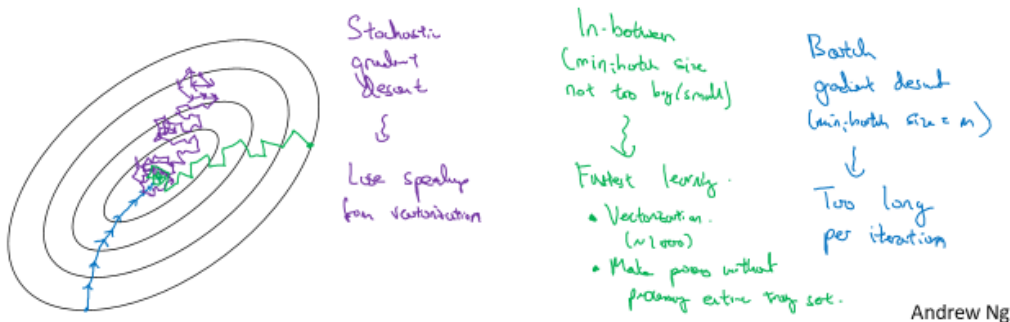- On mini batch gradient descent though, if you plot progress on your cost function, then it may not decrease on every iteration. In particular, on every iteration you're processing some X{t}, Y{t} and so if you plot the cost function J{t}, which is computer using just X{t}, Y{t}. Then it's as if on every iteration you're training on a different training set or really training on a different mini batch. So you plot the cross function J, you're more likely to see something that looks like this. It should trend downwards, but it's also going to be a little bit noisier.



Andrew Ng

In practice, somewhere between 1 and m (sample size)

If these are the contours of the cost function you're trying to minimize so your minimum is there. Then batch gradient descent might start somewhere and be able to take relatively low noise, relatively large steps. And you could just keep matching to the minimum. In contrast with stochastic gradient descent If you start somewhere let's pick a different starting point. Then on every iteration you're taking gradient descent with just a single strain example so most of the time you hit two at the global minimum. But sometimes you hit in the wrong direction if that one example happens to point you in a bad direction. So stochastic gradient descent can be extremely noisy. And on average, it'll take you in a good direction, but sometimes it'll head in the wrong direction as well. As stochastic gradient descent won't ever converge, it'll always just kind of oscillate and wander around the region of the minimum. But it won't ever just head to the minimum and stay there. In practice, the mini-batch size you use will be somewhere in between.

Choosing your mini-batch size

- If small training set (m<2000):  Use batch gradient descent
- If bigger training set
  Typical mini-batch size
  -> 64, 128, 256, 512

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

In practice of course the mini batch size is another hyper parameter that you might do a quick search over to try to figure out which one is most sufficient of reducing the cost function j. So what i would do is just try several different values. Try a few different powers of two and then see if you can pick one that makes your gradient descent optimization algorithm as efficient as possible. But it turns out there're even more efficient algorithms than gradient descent or mini-batch gradient descent. Let's start talking about them in the next few videos.

**2.2Exponentially weighted averages**

## Temperature in London

$\theta_1 = 40°F$ $\quad$ 4°C $\leftarrow$
$\theta_2 = 49°F$ $\quad$ 9°C
$\theta_3 = 45°F$ $\quad$ ⋮

⋮

$\theta_{180} = 60°F$ $\quad$ 15°C
$\theta_{181} = 56°F$ $\quad$ ⋮

⋮

$V_0 = 0$
$V_1 = 0.9\, V_0 + 0.1\, \theta_1$
$V_2 = 0.9\, V_1 + 0.1\, \theta_2$
$V_3 = 0.9\, V_2 + 0.1\, \theta_3$
⋮
$V_t = 0.9\, V_{t-1} + 0.1\, \theta_t$

## Exponentially weighted $\overset{moving}{averages}$

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

$\beta = 0.9$ $\quad : \quad$ ≈ 10 days' temperature
$\beta = 0.98$ $\quad : \quad$ ≈ 50 days
$\beta = 0.5$ $\quad : \quad$ ≈ 2 days

$V_t$ as approximate average over
$\rightarrow \approx \frac{1}{1-\beta}$ days' temperature.

$\frac{1}{1-0.98} = 50$

15

# Exponentially weighted averages

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

$\beta = 0.9 \qquad 0.98 \qquad 0.5$



temperature

days

# Exponentially weighted averages

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

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

…

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

$$= 0.1\,\theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1\,(0.9)^2\,\theta_{98} + 0.1\,(0.9)^3\,\theta_{97} + 0.1\,(0.9)^4\,\theta_{96}$$
$$+ \ldots$$

$\theta_t$

$t \sim \frac{1}{1-\beta}$

$V_{100}$   $10$   $\varepsilon = 1-\beta$

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

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

$(1-\varepsilon)^{1/\varepsilon} = \frac{1}{e}$    $0.98\,?$

$0.9 \quad \varepsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}$

# Implementing exponentially weighted averages

$$v_0 = 0$$
$$v_1 = \beta v_0 + (1-\beta)\theta_1$$
$$v_2 = \beta v_1 + (1-\beta)\theta_2$$
$$v_3 = \beta v_2 + (1-\beta)\theta_3$$

…

$$V_\theta := 0$$
$$V_\theta := \beta v + (1-\beta)\theta_1$$
$$V_\theta := \beta v + (1-\beta)\theta_2$$
$$\vdots$$

$\rightarrow V_\theta = 0$

Repeat {

    Get next $\theta_t$

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

}

16

**2.3 Bias correction in exponentially weighted averages**

# Bias correction



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

$$V_0 = 0$$

$$V_1 = \cancel{0.98 V_0} + 0.02\,\Theta_1$$

$$V_2 = 0.98\,V_1 + 0.02\,\Theta_2$$

$$= 0.98 \times 0.02 \times \Theta_1 + 0.02\,\Theta_2$$

$$= 0.0196\,\Theta_1 + 0.02\,\Theta_2$$

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

$$t=2: \quad 1-\beta^t = 1-(0.98)^2 = 0.0396$$

$$\frac{V_2}{0.0396} = \frac{0.0196\,\Theta_1 + 0.02\,\Theta_2}{0.0396}$$

$\beta = 0.98$

Andrew Ng

**2.4 Gradient descent with momentum**

# Gradient descent example



↑ slower learning
↔ faster learning

Momentum:
On iteration $t$:
  Compute $dW$, $db$ on current mini-batch.
  $V_{dW} = \beta V_{dW} + (1-\beta)dW$
  $V_{db} = \beta V_{db} + (1-\beta)db$
  $W = W - \alpha V_{dW}$, $b = b - \alpha V_{db}$

friction ⟵ ↗ velocity   ↗ acceleration

"$V_\theta = \beta V_{\theta\tau} + (1-\beta)\theta_t$"

deeplearning.ai

Andrew Ng

There's an algorithm called momentum, or gradient descent with momentum that almost always works faster than the standard gradient descent algorithm. In one sentence, the basic idea is to compute an exponentially weighted average of your gradients, and then use that gradient to update your weights instead. As an example let's say that you're trying to optimize a cost function which has contours like

this. So the red dot denotes the position of the minimum. Maybe you start gradient descent here and if you take one iteration of gradient descent either or descent maybe end up -heading there. But now you're on the other side of this ellipse, and if you take another step of gradient descent maybe you end up doing that. And then another step, another step, and so on. And you see that gradient descents will sort of take a lot of steps, right? Just slowly oscillate toward the minimum. And this up and down oscillations slows down gradient descent and prevents you from using a much larger learning rate. In particular, if you were to use a much larger learning rate you might end up over shooting and end up diverging like so. And so the need to prevent the oscillations from getting too big forces you to use a learning rate that's not itself too large. Another way of viewing this problem is that on the vertical axis you want your learning to be a bit slower, because you don't want those oscillations. But on the horizontal axis, you want faster learning.

If you average out these gradients, you find that the oscillations in the vertical direction will tend to average out to something closer to zero. So, in the vertical direction, where you want to slow things down, this will average out positive and negative numbers, so the average will be close to zero. Whereas, on the horizontal direction, all the derivatives are pointing to the right of the horizontal direction, so the average in the horizontal direction will still be pretty big. So that's why with this algorithm, with a few iterations you find that the gradient descent with momentum ends up eventually just taking steps that are much smaller oscillations in the vertical direction, but are more directed to just moving quickly in the horizontal direction. And so this allows your algorithm to take a more straightforward path, or to damp out the oscillations in this path to the minimum.

# Implementation details

$V_{dw} = 0$ , $V_{db} = 0$

On iteration $t$:

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

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

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

$W = W - \alpha v_{dW}, \; b = b - \alpha v_{db}$

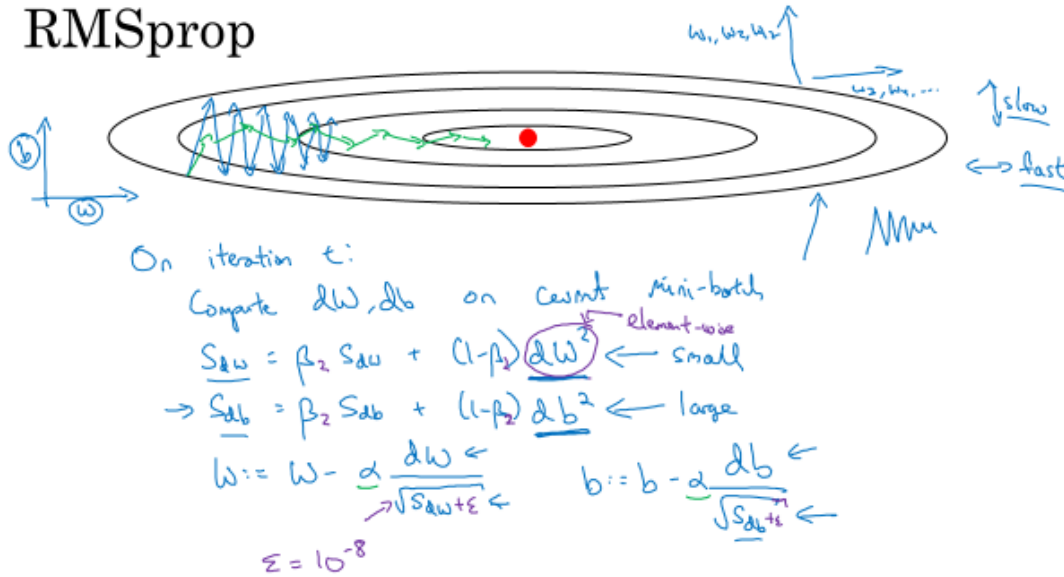Hyperparameters: $\alpha, \beta$      $\beta = 0.9$

average over last $\times 10$ gradients

Andrew Ng

## 2.5 RMSprop

RMSprop: root mean squre prop, which also can speed up gradient descent (like momentum)

# RMSprop



On iteration $t$:

Compute $dW, db$ on current mini-batch

$S_{dW} = \beta_2 S_{dW} + (1-\beta_2)(dW)^2 \leftarrow$ small (element-wise)

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

$W := W - \alpha \dfrac{dW}{\sqrt{S_{dW}} + \varepsilon} \leftarrow \qquad b := b - \alpha \dfrac{db}{\sqrt{S_{db}} + \varepsilon} \leftarrow$

$\varepsilon = 10^{-8}$

That your updates in the, Vertical direction and then horizontal direction you can keep going. And one effect of this is also that you can therefore use a larger learning rate alpha, and get faster learning without diverging in the vertical direction.

So that's RMSprop, and similar to momentum, has the effects of damping out the oscillations in gradient descent, in mini-batch gradient descent. And allowing you to maybe use a larger learning rate alpha. And certainly speeding up the learning speed of your algorithm. One fun fact about RMSprop, it was actually first proposed not in an academic research paper, but in a Coursera course that Jeff Hinton had taught on Coursera many years ago. I guess Coursera wasn't intended to be a platform for dissemination of novel academic research, but it worked out pretty well in that case. And was really from the Coursera course that RMSprop started to become widely known and it really took off.

**2.6 Adam optimization algorithm**

# Adam optimization algorithm

$V_{dW} = 0, \; S_{dW} = 0. \quad V_{db} = 0, \; S_{db} = 0$

On iteration $t$:

Compute $dW, db$ using current mini-batch

$V_{dW} = \beta_1 V_{dW} + (1-\beta_1) dW \quad , \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \leftarrow$ "momentum" $\beta_1$

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

yhat = np.array([.9, 0.2, 0.1, .4, .9])

$V_{dW}^{corrected} = V_{dW}/(1-\beta_1^t) \quad , \quad V_{db}^{corrected} = V_{db}/(1-\beta_1^t)$

$S_{dW}^{corrected} = S_{dW}/(1-\beta_2^t) \;,\; S_{db}^{corrected} = S_{db}/(1-\beta_2^t)$

$W := W - \alpha \dfrac{V_{dW}^{corrected}}{\sqrt{S_{dW}^{corrected}} + \varepsilon} \qquad b := b - \alpha \dfrac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected}} + \varepsilon}$

# Hyperparameters choice:

$$\rightarrow \alpha : \text{needs to be tune}$$
$$\rightarrow \beta_1 : 0.9 \quad \longrightarrow (dw)$$
$$\rightarrow \beta_2 : 0.999 \quad \longrightarrow (dw^2)$$
$$\rightarrow \varepsilon : 10^{-8}$$

Adam: Adaptive Moment Estimation

When implementing Adam, what people usually do is just use the default value. So, ß1 and ß2 as well as epsilon. I don't think anyone ever really tunes Epsilon. And then, try a range of values of Alpha to see what works best. You could also tune ß1 and ß2 but it's not done that often among the practitioners I know. ß1 is computing the mean of the derivatives. This is called the first moment. And ß2 is used to compute exponentially weighted average of the ²s and that's called the second moment. So that gives rise to the name adaptive moment estimation. But everyone just calls it the Adam authorization algorithm.

## 2.7 Learning rate decay

One of the things that might help speed up your learning algorithm, is to slowly reduce your learning rate over time. We call this learning rate decay.

. Let's start with an example of why you might want to implement learning rate decay. Suppose you're implementing mini-batch gradient descent, with a reasonably small mini-batch. Maybe a mini-batch has just 64, 128 examples. Then as you iterate, your steps will be a little bit noisy. And it will tend towards this minimum over here, but it won't exactly converge. But your algorithm might just end up wandering around, and never really converge, because you're using some fixed value for alpha. And there's just some noise in your different mini-batches. But if you were to slowly reduce your learning rate alpha, then during the initial phases, while your learning rate alpha is still large, you can still have relatively fast learning. But then as alpha gets smaller, your steps you take will be slower and smaller. And so you end up oscillating in a tighter region around this minimum, rather than wandering far away, even as training goes on and on.

So the intuition behind slowly reducing alpha, is that maybe during the initial steps of learning, you could afford to take much bigger steps. But then as learning approaches converges, then having a slower learning rate allows you to take smaller steps. So here's how you can implement learning rate decay. Recall that one epoch is one pass,

Learning rate decay



Slowly reduce $\alpha$

# Learning rate decay

1 epoch = 1 pass through data.

$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} \alpha_0 \quad \leftarrow$$

| Epoch | $\alpha$ |
|---|---|
| 1 | 0.1 |
| 2 | 0.67 |
| 3 | 0.5 |
| 4 | 0.4 |
| ⋮ | ⋮ |

$\boxed{X^{\{1\}}} \boxed{X^{\{2\}}} \quad \cdots$ → epoch 1
→ epoch 2

$\alpha_0 = 0.2$
decay-rate = 1

# Other learning rate decay methods

formula
$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \qquad - \text{exponentially decay.}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \qquad \text{or} \qquad \frac{k}{\sqrt{t}} \cdot \alpha_0$$

discrete staircase

Manual decay.

21

## 2.8 The problem of local optima



Local optima in neural networks

Andrew Ng

Local optima normally is not a problem.

In the early days of deep learning, people used to worry a lot about the optimization algorithm getting stuck in bad local optima. But as this theory of deep learning has advanced, our understanding of local optima is also changing. It turns out that if you are plotting a figure like this in two dimensions, then it's easy to create plots like this with a lot of different local optima. And these very low dimensional plots used to guide their intuition. But this intuition isn't actually correct. It turns out if you create a neural network, most points of zero gradients are not local optima like points like this. Instead most points of zero gradient in a cost function are saddle points. So, that's a point where the zero gradient, again, just is maybe W1, W2, and the height is the value of the cost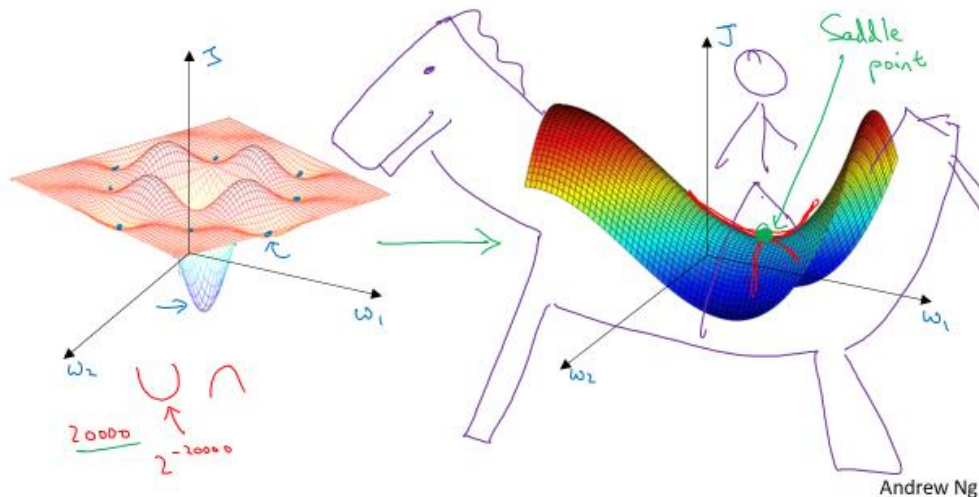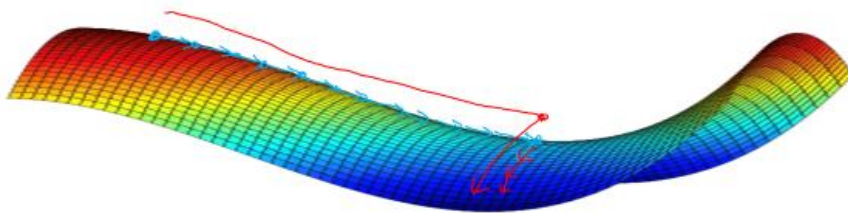 function J. But informally, a function of very high dimensional space, if the gradient is zero, then in each direction it can either be a convex light function or a concave light function. And if you are in, say, a 20,000 dimensional space, then for it to be a local optima, all 20,000 directions need to look like this. And so the chance of that happening is maybe very small, maybe two to the minus 20,000. Instead you're much more likely to get some directions where the curve bends up like so, as well as some directions where the curve function is bending down rather than have them all bend upwards. So that's why in very high-dimensional spaces you're actually much more likely to run into a saddle point like that shown on the right, then the local optimum. As for why the surface is called a saddle point, if you can picture, maybe this is a sort of saddle you put on a horse, right?

And so, one of the lessons we learned in history of deep learning is that a lot of our intuitions about low-dimensional spaces, like what you can plot on the left, they really don't transfer to the very high-dimensional spaces that any other algorithms are operating over. Because if you have 20,000 parameters, then J as your function over 20,000 dimensional vector, then you're much more likely to see saddle points than local optimum. If local optima aren't a problem, then what is a problem? It turns out that plateaus can really slow down learning and a plateau is a region where the derivative is close to

zero for a long time. So if you're here, then gradient descents will move down the surface, and because the gradient is zero or near zero, the surface is quite flat. You can actually take a very long time, you know, to slowly find your way to maybe this point on the plateau. And then because of a random perturbation of left or right, maybe then finally I'm going to search pen colors for clarity. Your algorithm can then find its way off the plateau. Let it take this very long slope off before it's found its way here and they could get off this plateau.

So the takeaways from this video are, first, you're actually pretty unlikely to get stuck in bad local optima so long as you're training a reasonably large neural network, save a lot of parameters, and the cost function J is defined over a relatively high dimensional space. But second, that plateaus are a problem and you can actually make learning pretty slow. And this is where algorithms like momentum or RmsProp or Adam can really help your learning algorithm as well. And these are scenarios where more sophisticated observation algorithms, such as Adam, can actually speed up the rate at which you could move down the plateau and then get off the plateau. So because your network is solving optimizations problems over such high dimensional spaces, to be honest, I don't think anyone has great intuitions about what these spaces really look like, and our understanding of them is still evolving. But I hope this gives you some better intuition about the challenges that the optimization algorithms may face.

## Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

**Week 3 Hyperparameter tuning, Batch Normalization and Programming Frameworks**

**3.1 Tuning process**

# Hyperparameters



The learning rate, alpha, is the most important hyper-parameter to tune.

Momentum beta is importance, normally the default value 0.9 is fine

Hyper-parameters beta1, beta2 for Adam are important.

- Red: First important
- Yellow: Second important
- Purple: Third important

# Try random values: Don't use a grid



Coarse to fine

**3.2Using an appropriate scale to pick hyper-parameters**

## Picking hyperparameters at random

$\rightarrow n^{[l]} = 50, \ldots, 100$

$\underset{50 \qquad\qquad\qquad\qquad 100}{[\;*\;\; *\!*\!*\;\; \times\!\times\;\; \times\;\times\;\; \times\;]}$

$\rightarrow$ #layers $\quad L: \quad 2-4$

$\qquad\qquad 2, 3, 4$

## Appropriate scale for hyperparameters

$\alpha = 0.0001 \ldots, 1$

$\underset{0.0001 \qquad\qquad\qquad\qquad\qquad 1}{[\;\times\;\;\times\;\times\;\times\;\times\;\;\times\times\;\;\times\times\times\;\;\times\;\;\times\times\;]}$

$\underset{0.1}{\smile} \cdots \; - \; - \; - \; - \; - \; - \cdot \cdot 1$

$\underset{0.0001}{[\;\times\;\;\times\;|\times\times\;\;\times?\;|\;\times\times\times\;\times\;|\;\times\times\;]}$

$\boxed{0.0001} \quad \underline{0.001} \quad \underline{0.01} \quad \underline{0.1} \quad \boxed{1}$

$\quad 10^{-4} \qquad\qquad\qquad\qquad\qquad\qquad 10^{0} \quad 10^{6} \quad b = \log_{10} 1$

$10^{a} \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = 0$

$a = \log_{10} 0.0001 \quad r = -4 * np.random.rand() \quad \leftarrow \quad r \in [-4, 0]$

$= -4 \qquad \alpha = 10^{r} \qquad\qquad\qquad \leftarrow \quad 10^{-4} \ldots 10^{0}$
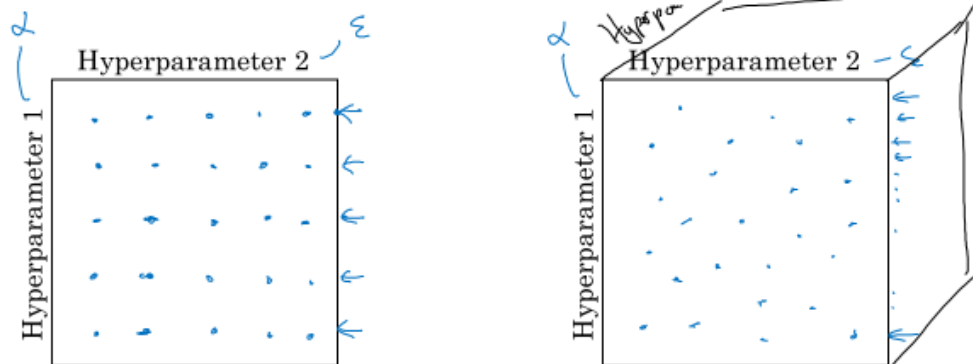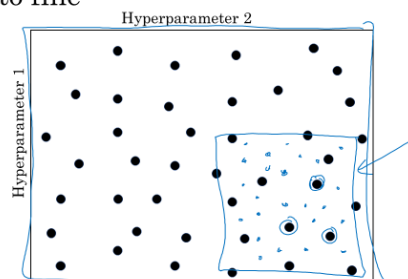
$\underline{10^{a} \ldots 10^{b}} \qquad \underline{\underset{[-4, 0]}{r \in [a, b]}} \qquad \underline{\alpha = 10^{r}}$

## Hyperparameters for exponentially weighted averages

$\beta = 0.9 \quad \ldots \quad 0.999$

$\qquad \downarrow \qquad\qquad\qquad \downarrow$

$\qquad 10 \qquad\qquad\qquad 1000$

$\underset{0.9 \qquad\qquad\qquad 0.999}{[\;\times\!\times\;\times\!\times\!\times\!\times\;]} \leftarrow$

$\underset{0.9 \qquad 0.99 \qquad 0.999}{[\underline{\qquad}\,\underline{\qquad}\,]}$

$1 - \beta = 0.1 \quad \ldots \quad 0.001$

$\underset{0.1 \qquad 0.01 \qquad 0.001}{[\underline{\qquad}\,\underline{\qquad}\,]}$

$\qquad\qquad\qquad\qquad\qquad\qquad 10^{-1} \qquad\qquad 10^{-3}$

$\beta: \; 0.900 \rightarrow 0.9005 \; \Big\} \sim 10$

$\beta: \; 0.999 \rightarrow 0.9995 \qquad\qquad r \in [-3, -1]$

$\quad \sim 1000 \qquad \sim 2000 \qquad \frac{1}{1-\beta} \qquad 1 - \beta = 10^{r}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \beta = 1 - 10^{r}$

### 3.3 Hyperparameters tuning in practice Pandas vs. Caviar

## Re-test hyperparameters occasionally



- NLP, Vision, Speech, Ads, logistics, ….

- Intuitions do get stale. Re-evaluate occasionally.

Babysitting one model

Training many models in parallel

Panda

Caviar

Andrew Ng

### 3.4 Normalizing activations in a network

Batch normalization

## Normalizing inputs to speed up learning



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

$$X = X - \mu \quad \text{element-wise}$$

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

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

Can we normalize $a^{[2]}$ so to train $w^{[3]}, b^{[3]}$ faster

Normalize $z^{[2]}$

Andrew Ng

We actually normalize the value before activation function (Z), not after (a).

# Implementing Batch Norm

Given some intermediate values in NN $\quad z^{(1)}, \ldots, z^{(m)}$

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

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

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

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

Use $\tilde{z}^{[l](i)}$ instead of $z^{[l](i)}$.

If $\gamma = \sqrt{\sigma^2 + \varepsilon}$

$\beta = \mu$

then $\tilde{z}^{(i)} = z^{(i)}$

learnable parameters of model.

$X \leftarrow$
$z^{(i)} \leftarrow$

One difference between the training input and these hidden unit values is you might not want your hidden unit values be forced to have mean 0 and variance 1. For example, if you have a sigmoid activation fu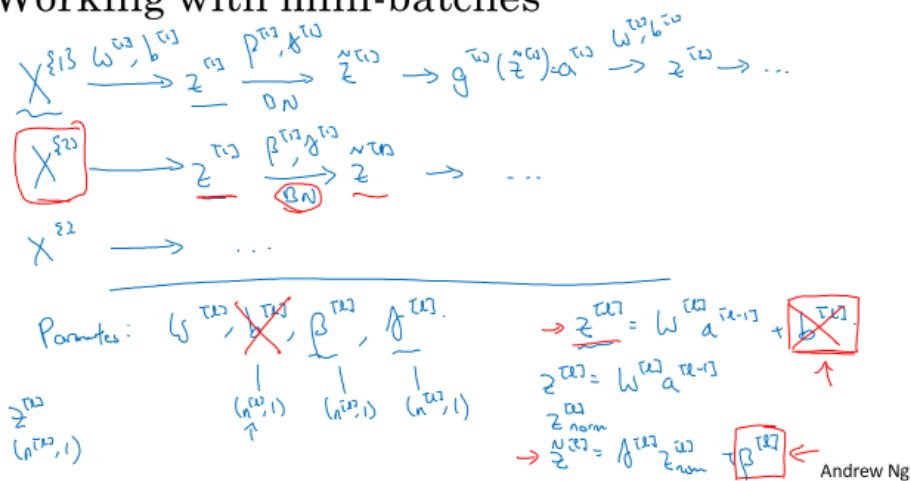nction, you don't want your values to always be clustered here. You might want them to have a larger variance or have a mean that's different than 0, in order to better take advantage of the nonlinearity of the sigmoid function rather than have all your values be in just this linear regime. So that's why with the parameters gamma and beta, you can now make sure that your zi values have the range of values that you want. But what it does really is it then shows that your hidden units have standardized mean and variance, where the mean and variance are controlled by two explicit parameters gamma and beta which the learning algorithm can set to whatever it wants. So what it really does is it normalizes in mean and variance of these hidden unit values, really the zis, to have some fixed mean and variance. And that mean and variance could be 0 and 1, or it could be some other value, and it's controlled by these parameters gamma and beta.

### 3.5 Fitting Batch Norm into a neural network

$$X \xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \to a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{w^{[2]}, b^{[2]}} z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \to a^{[2]} \to$$

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \ldots, w^{[l]}, b^{[l]},$
$\to \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \ldots, \beta^{[l]}, \gamma^{[l]}$

$\to \beta$

$d\beta^{[l]}$

$\beta = \beta - \alpha \, d\beta^{[l]}$

tf.nn.batch-normalization $\leftarrow$

# Working with mini-batches

$$X^{\{1\}} \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \to g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \to \dots$$

$$X^{\{2\}} \longrightarrow Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \to \dots$$

$$X^{\{2\}} \longrightarrow \dots$$

Paramtes: $W^{[1]}$, $b^{[1]}$, $\beta^{[1]}$, $\gamma^{[1]}$.

$Z^{[l]}$
$(n^{[l]}, 1)$

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

$$\to Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$
$$Z^{[l]} = W^{[l]} a^{[l-1]}$$
$$Z^{[l]}_{norm}$$
$$\to \tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]}_{norm} + \beta^{[l]} \leftarrow$$

Andrew Ng

# Implementing gradient descent
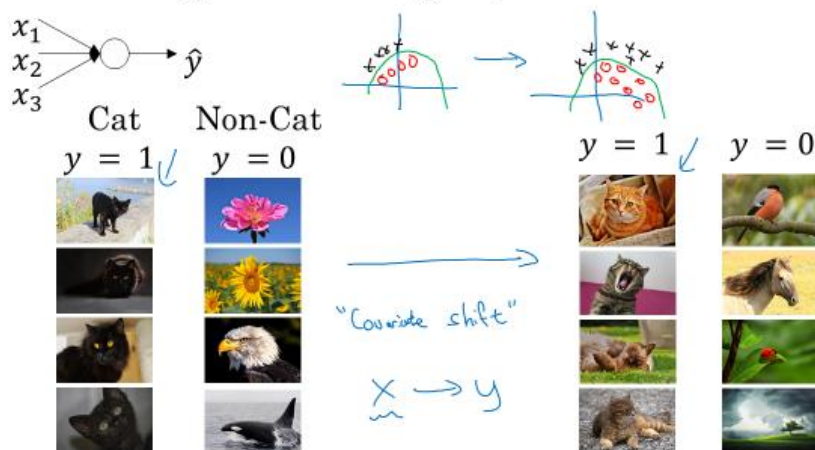
for $t = 1 \dots$ num MiniBatches
    Compute forward prop on $X^{\{t\}}$.
       In each hidden layer, use BN to repca $Z^{[l]}$ with $\tilde{Z}^{[l]}$.
   Use backprop to compte $dW^{[l]}$, $db^{[l]}$, $d\beta^{[l]}$, $d\gamma^{[l]}$
   Update parates $W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$
$$\left. \begin{array}{l} \beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]} \\ \gamma^{[l]} := \dots \end{array} \right\} \leftarrow$$

Works w/ momentm, RMSprop, Adam.

**3.6Why does Batch Norm work**

# Learning on shifting input distribution

$x_1$
$x_2$ $\to \hat{y}$
$x_3$

Cat    Non-Cat
$y = 1$    $y = 0$

"Covriate shift"

$$x \to y$$

$y = 1$    $y = 0$

Why does batch norm work? Here's one reason, you've seen how normalizing the input features, the X's, to mean zero and variance one, how that can speed up learning. So rather than having some features that range from zero to one, and some from one to a 1,000, by normalizing all the features, input features X, to take on a similar range of values that can speed up learning. So, one intuition behind why batch norm works is, this is doing a similar thing, but further values in your hidden units and not just for your input there.

This idea of your data distribution changing goes by the somewhat fancy name, covariate shift. And the idea is that, if you've learned some X to Y mapping, if the distribution of X changes, then you might need to retrain your learning algorithm. And this is true even if the function, the ground true function, mapping from X to Y, remains unchanged, which it is in this example, because the ground true function is, is this picture a cat or not. And the need to retain your function becomes even more acute or it becomes even worse if the ground true function shifts as well.

# Why this is a problem with neural networks?



# Batch Norm as regularization

- 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.

**3.7 Batch Norm at test time**

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

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

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

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

$\mu, \sigma^2$: estimate using exponentially weighted average (across mini-batches).

$X^{\{1\}}, X^{\{1\}}, X^{\{3\}}, \dots$

$\mu^{\{1\}[l]} \quad \mu^{\{2\}[l]} \quad \mu^{\{3\}[l]} \longrightarrow \mu$

$\theta_1 \quad \theta_2 \quad \theta_3 \quad \sigma^2$

$\sigma^{2\{1\}[l]} \quad \sigma^{2\{2\}[l]} \quad \dots$

$z_{norm} = \frac{z - \mu}{\sqrt{\sigma^2 + \varepsilon}} \quad\quad \tilde{z} = \gamma z_{norm} + \beta$

So the takeaway from this is that during training time mu and sigma squared are computed on an entire mini batch of say 64 engine, 28 or some number of examples. But that test time, you might need to process a single example at a time. So, the way to do that is to estimate mu and sigma squared from your training set and there are many ways to do that. You could in theory run your whole training set through your final network to get mu and sigma squared. But in practice, what people usually do is implement and exponentially weighted average where you just keep track of the mu and sigma squared values you're seeing during training and use and exponentially the weighted average, also sometimes called the running average, to just get a rough estimate of mu and sigma squared and then you use those values of mu and sigma squared that test time to do the scale and you need the head and unit values Z. In practice, this process is pretty robust to the exact way you used to estimate mu and sigma squared. So, I wouldn't worry too much about exactly how you do this and if you're using a deep learning framework, they'll usually have some default way to estimate the mu and sigma squared that should work reasonably well as well. But in practice, any reasonable way to estimate the mean and variance of your head and unit values Z should work fine at test.
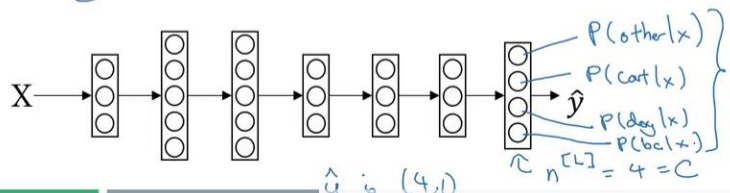
**3.8 Softmax Regression**

There's a generalization of logistic regression called Softmax regression. The less you make predictions where you're trying to recognize one of C or one of multiple classes, rather than just recognize two classes

# Recognizing cats, dogs, and baby chicks, other 0
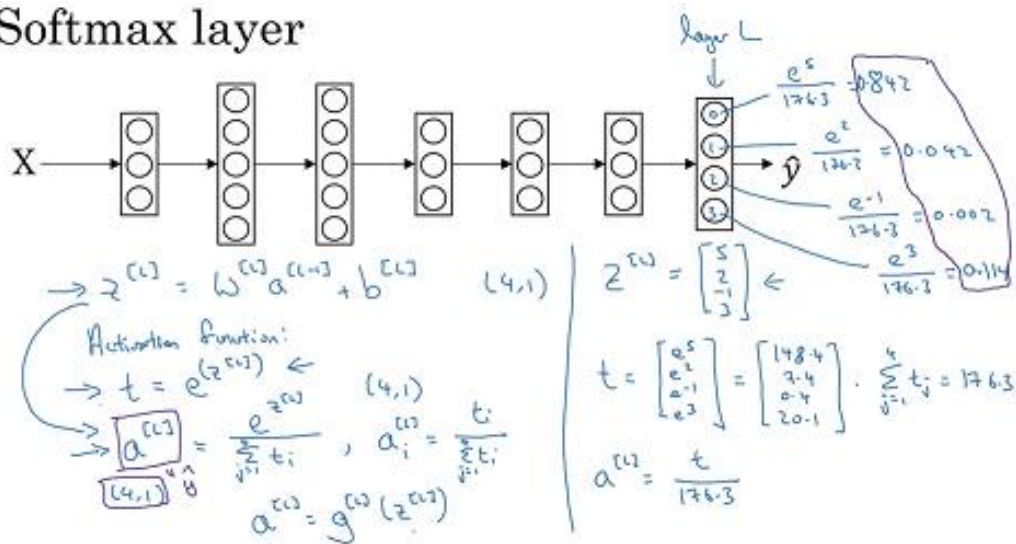


$$C = \#classes = 4 \qquad (0, \dots, 3)$$



$P(other|x)$
$P(cat|x)$
$P(dog|x)$
$P(bc|x)$

$n^{[L]} = 4 = C$

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

# Softmax layer

layer L



$\frac{e^5}{176.3} = 0.842$

$\frac{e^2}{176.3} = 0.042$

$\frac{e^{-1}}{176.3} = 0.002$

$\frac{e^3}{176.3} = 0.114$

$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]} \qquad (4,1)$

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

Activation function:

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

$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^{4} t_i}, \qquad a_i^{[L]} = \frac{t_i}{\sum t_i}$

$(4,1) \, \hat{y}$

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

$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \sum_{j=1}^{4} t_j = 176.3$

$a^{[L]} = \frac{t}{176.3}$

# Softmax examples

31

# Understanding softmax

$$(4,1)$$

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \qquad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \qquad\qquad C = 4 \qquad\qquad g^{[L]}(\cdot)$$

"soft max"

"hard max"

$$a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5/(e^5+e^2+e^{-1}+e^3) \\ e^2/(e^5+e^2+e^{-1}+e^3) \\ e^{-1}/(e^5+e^2+e^{-1}+e^3) \\ e^3/(e^5+e^2+e^{-1}+e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \qquad\qquad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$
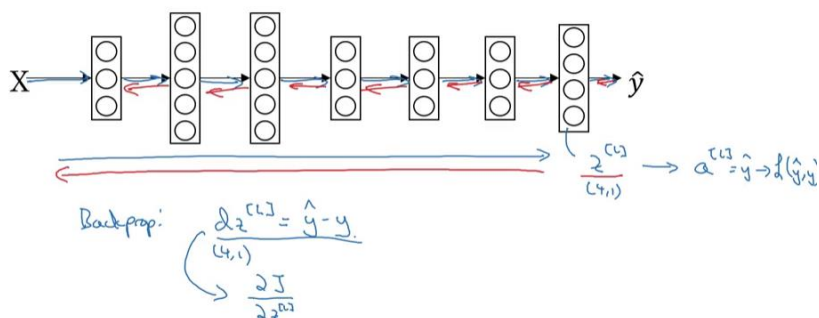
Softmax regression generalizes logistic regression to $C$ classes.

If $\underline{C=2}$, softmax reduces to logistic regression. $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$

# Loss function

$$(4,1)$$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \leftarrow \text{cat} \qquad y_2 = 1 \qquad\qquad (4,1)$$

$$y_1 = y_3 = y_4 = 0$$

$$a^{[L](i)} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \qquad\qquad C = 4$$

$$\mathcal{L}(\hat{y}, y) = -\sum_{j=1}^{4} y_j \log \hat{y}_j \qquad\qquad J(w^{[1]}, b^{[1]}, \dots) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

small

$$-y_2 \log \hat{y}_2 = -\log \hat{y}_2. \qquad\qquad \text{Make } \hat{y}_2 \text{ big.}$$

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

$$= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \cdots \qquad\qquad = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \cdots$$

$$(4, m) \qquad\qquad\qquad (4, m)$$

Andrew Ng

# Gradient descent with softmax



$$X \longrightarrow \cdots \longrightarrow \hat{y}$$

$$\underset{(4,1)}{z^{[L]}} \longrightarrow a^{[L]} = \hat{y} \longrightarrow \mathcal{L}(\hat{y}, y)$$

Backprop: $\underset{(4,1)}{dz^{[L]}} = \hat{y} - y.$

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

32

### 3.10 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)

### 3.11 TensorFlow

## Motivating problem

$$J(w) = \boxed{w^2 - 10w + 25}$$
(cost)

$$\llcorner (w-5)^2$$

$$w = 5$$

$$J(w,b)$$
↑ ↑

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

In [5]:  coefficients = np.array([[1.], [-10.], [25.]])

         w = tf.Variable(0,dtype=tf.float32)
         x = tf.placeholder(tf.float32, [3,1])
         #cost = tf.add(tf.add(w**2,tf.multiply(-10.,w)),25)
         #cost = w**2 - 10*w + 25
         cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]
         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 [9]:  session.run(train, feed_dict={x:coefficients})
         print(session.run(w))

         0.1

In [10]: for i in range(1000):
             session.run(train, feed_dict={x:coefficients})
         print(session.run(w))

         4.99999
```

# 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))
```
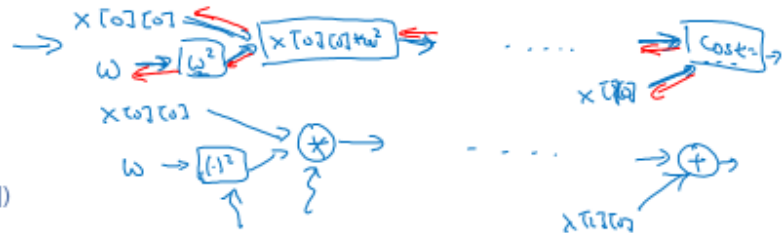
Andrew Ng

34