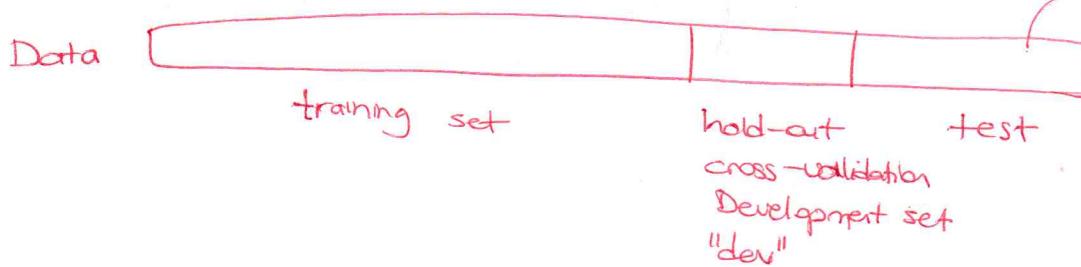


Train / Dev / Test sets

1

Applied ML is a highly iterative process

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



The goal of the test set is to give you a unbiased estimate of the performance of your final network.

\* The goal is to see which of many different models performs best on your dev set ✓

### Mismatched train/test distribution

#### Training set:

Cat pictures from websites  
 ↴  
 high resolution  
 very professional  
 very nicely framed pictures

#### Dev/test sets:

Cat pictures from users using your app  
 ↴  
 blurrier  
 low resolution

⇒ So these two distributions of data may be different

#### ! Rule :

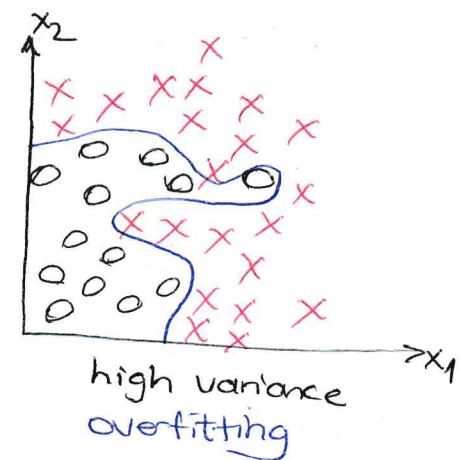
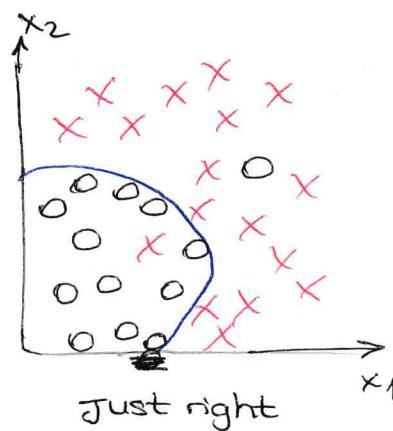
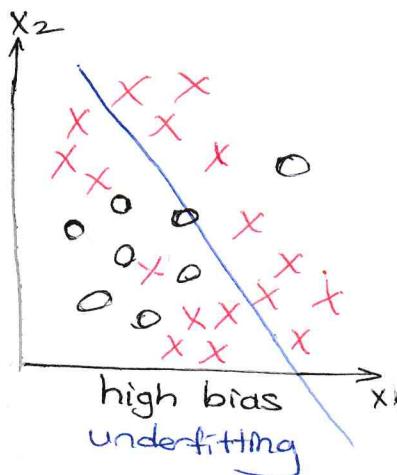
Make sure dev and test sets ~~come from~~ the same distribution

→ Not having a test set might be okay  
 (Only dev set)

train / development set would be more correct terminology

So: Having set up a train, dev and test set will allow you to integrate more quickly.

It will also allow you to more efficiently measure the bias and variance of your algorithm. So ways can be selected more efficiently to improve algorithms.

Bias / Variance

For Example : CAT CLASSIFICATION

Train set error : 1%

Dev set error : 11%

high variance

15%

16%

high bias

15%

30%

high bias

& high variance  
WORST

0.5%

1%

low bias

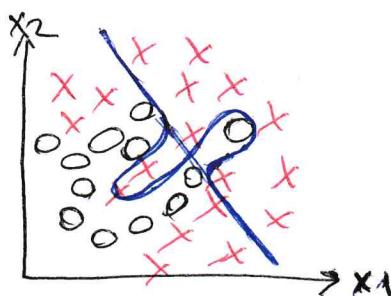
low variance

Human : ≈ 0%

Optimal (Base) error : ≈ 0%

By looking at training set error, you can get a sense of how well you are fitting. So that it tells you have a bias problem.

With dev set error, that should give you a sense of how bad is the variance problem.

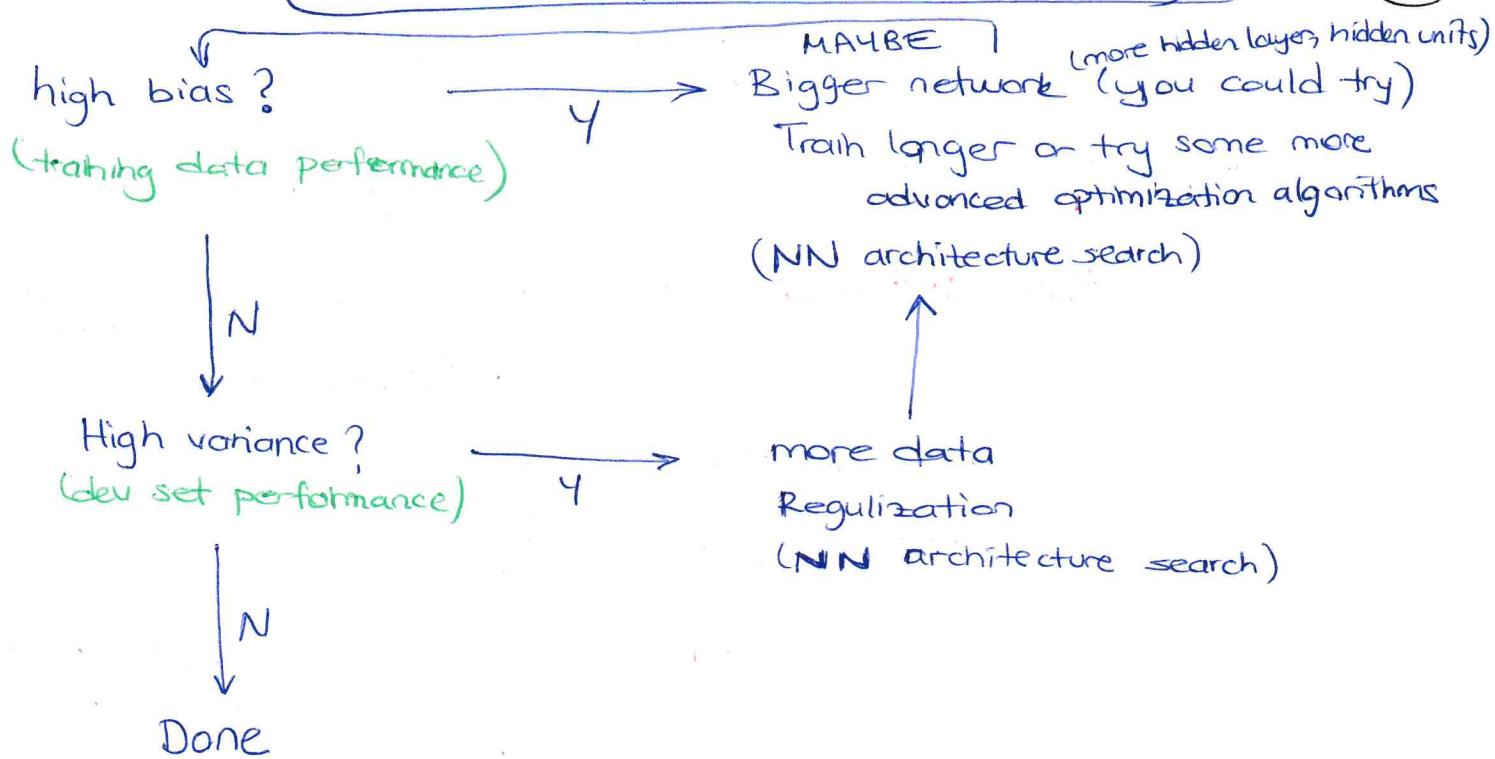


high ~~bias~~  
high variance

(too general)  
(too much flexibility)

Basic Recipe For Machine Learning

(3)

"bias variance tradeoff"

\* Training a bigger network almost never hurts. And the main cost of training a neural network that's too big is just computational time, so long as you're regularizing.

Regularization is very useful technique while reducing the variance. There is a little bit of a bias variance tradeoff when you're using regularization.

Regularization

Overfitting → high variance problem

Logistic Reg :

$$\min_{w,b} J(w,b) \quad w \in \mathbb{R}^n, b \in \mathbb{R}$$

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

↑ penalizes weight matrix from being too much

$$L_2 \text{ regularization } \|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$$

(most common type)

$$L_1 \text{ regularization } \frac{\lambda}{2m} \sum_{j=1}^n |w_j| = \frac{\lambda}{2m} \|w\|_1$$

$\lambda$  = regularization parameter

"lambd"

## Neural Network:

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

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

"Frobenius norm" means the sum of squares of elements of a matrix)

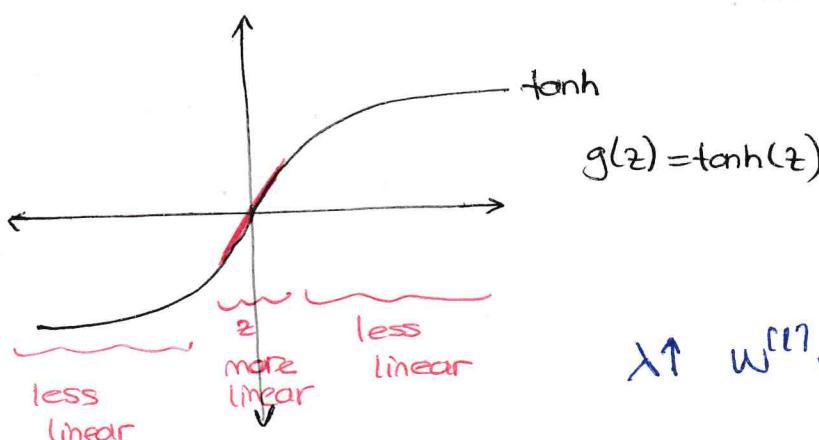
L2 regularization is sometimes also called "weight decay"

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

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

So you are multiplying the weight metrics by a number of slightly less than 1.

### Why Regularization reduces overfitting?



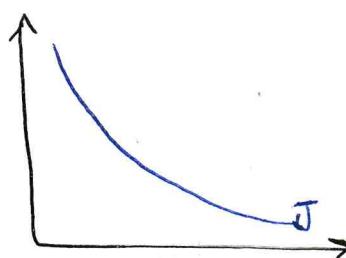
If  $\lambda$  is too high then  $w$  would be  $\approx 0$ . It would make the some hidden units smaller and gets the smaller neural network.

$$\lambda \uparrow \quad w^{(l)} \downarrow \quad z^{(l)} = w^{(l)} \cdot a^{(l-1)} + b^{(l)}$$

↓

will look like a simpler model like linear regression

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



Dropout RegularizationImplementing dropout ("Inverted Dropout")

Illustrate with layer=3 keep\_prob=0.8

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

$a3 = np.multiply(a3, d3)$  #  $a3 \neq d3$

$a3 /= keep\_prob$

most common  
dropout implementa  
~~tion~~

initially, 50 units  $\rightarrow$  10 units shut off

$$z^{(u)} = w^{(u)} \cdot a^{(3)} + b^{(u)}$$

This will be reduced by 0.20

Making Predictions at Test Time:

No drop out

~~If you want it at test time, it just add noise to your predictions.~~

Understanding DropoutWhy does drop-out work?

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

Dropout has a similar effect to L2 regularization

Shrink weights  
similar L2 reg.

More adaptive to  
the scale of different  
inputs.

It is also feasible to vary keep-prop by layer.

For the layer with a lot of parameters, keep-prop can be set smaller to apply a more powerful form of dropout.

Many of the first successful implementations of dropouts were to computer

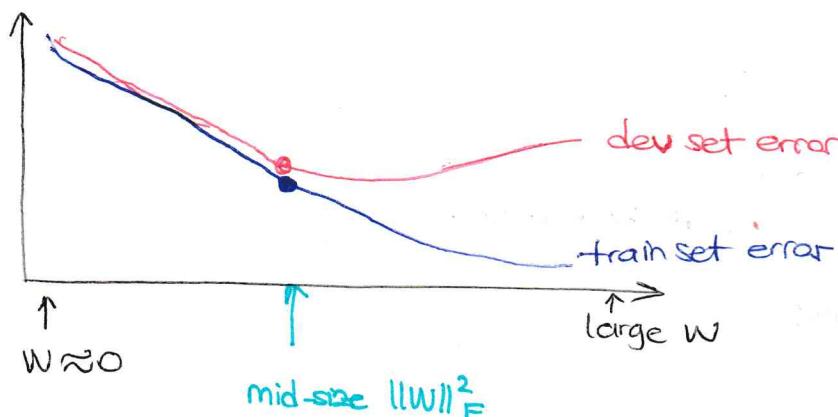
In computer vision input size is so big, inputting all these pixels that you almost never have enough data. So dropout is very frequently used.

Other Regularization MethodsData Augmentation:

Image of cat can be rotated or randomly zoomed in order to get more data.

This can be an inexpensive way to give your algorithm more data.

So data augmentation can be used as a regularization technique

Early Stopping:

Early stopping does let you get similar effect without needing to explicitly try lots of different values of lambda.

Orthogonalization

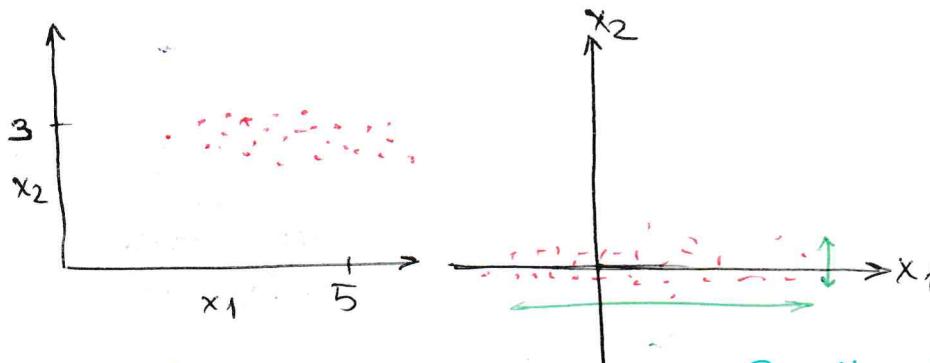
- Optimize the cost function J

Gradient Descent, momentum,

RMS prop, Adam ...

- Not overfit

Regularization, getting more data ...

Normalizing Inputs

1. Subtract mean:

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

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

2. Normalize variance:

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

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

element-wise

For train and test set, some steps should be used.

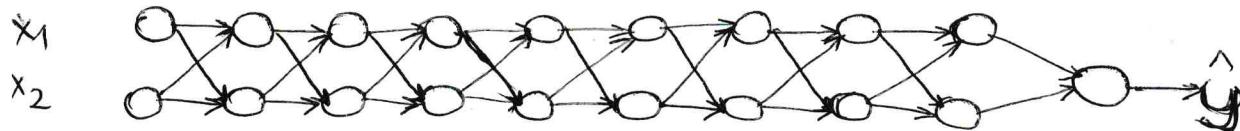
If you don't normalize, you should keep small the learning rate. But when inputs are normalized, learning-rate can be larger.

Normalizing makes your cost function J easier and faster to optimize.

All your features on a similar scale will usually help your Learning algorithm run faster.

Vanishing/exploding Gradients

One of the problems of training neural network especially very deep neural networks, is data vanishing and exploding gradients.



$$g(z) = z \quad b^{(l)} = 0$$

$$\hat{y} = W^{(l)} W^{(l-1)} \dots W^{(3)} [W^{(2)} W^{(1)} X]$$

$W^{(l)} > I \Rightarrow$  exploding

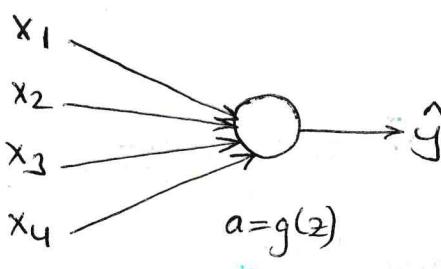
$W^{(l)} < I \Rightarrow$  vanishing

$$z^{(l)} = W^{(l)} X$$

$$a^{(l)} = g(z^{(l)}) = z^{(l)}$$

$$\underline{a^{(2)}} = g(z^{(2)}) = g(W^{(2)} a^{(1)})$$

If your activations or gradients increase or decrease exponentially as a function of L, then these values could get really big or really small. And this makes training difficult, especially if your gradients are exponentially smaller than 1 then gradient descent will take tiny little steps. It will take a long time for gradient descent to learn anything.

Weight Initialization for Deep Networks

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

large n  $\rightarrow$  smaller  $w_i$

$$\text{Var}(w_i) = \frac{1}{n} \rightarrow \text{for ReLU}$$

$$W^{(l)} = \text{np.random.rand(shape)} * \sqrt{\frac{2}{n^{(l-1)}}}$$

Other variants:

$$\tanh = \sqrt{\frac{1}{n^{(l-1)}}}$$

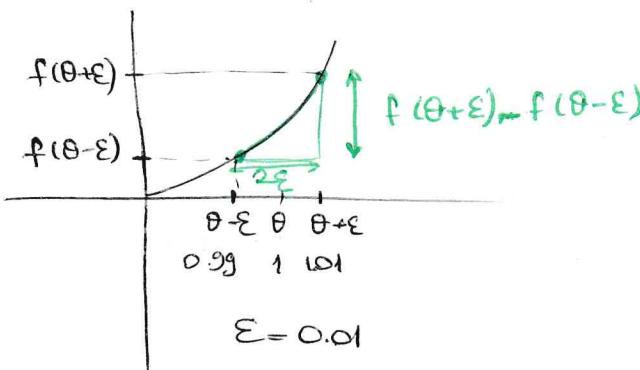
called "Xavier"

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

$\rightarrow$  by Yoshua Bengio and colleagues

Variance parameter could be another thing that you could tune of your hyperparameters

With these functions  $W^{(l)}$  won't be too much bigger than 1 and too much less than 1. So it doesn't explode or vanish too quickly.

Numerical Approximation Of Gradients

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

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

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

approximation error = 0.0001

Note:  $f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$

(prev. slide: 3.0301 error: 0.03)

$$\underbrace{\qquad}_{\text{more accurate}} \rightarrow O(\epsilon^2)$$

$$\begin{matrix} 0.01 \\ 0.0001 \end{matrix}$$

$$\lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \text{error: } O(\epsilon) \quad \begin{matrix} 0.01 \\ 0.001 \end{matrix}$$

less accurate

Gradient Checking

You can numerically verify whether or not a function  $g$ ,  $g(\theta)$  is a correct implementation of the  $f'(2)$

Take  $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$  and reshape into a big vector  $\theta$ .  
 $J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = J(\theta)$   
 Take  $dW^{(1)}, dB^{(1)}, \dots, dW^{(L)}, dB^{(L)}$  and reshape into a big vector  $d\theta$ .

It is used to find bugs in the implementation of backpropagation.

for each  $i$ :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1^{\epsilon}, \theta_2^{\epsilon}, \dots, \theta_i^{\epsilon}, \dots) - J(\theta_1^{\epsilon}, \theta_2^{\epsilon}, \dots, \theta_i^{-\epsilon}, \dots)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{\text{approx}} \approx d\theta.$$

Check

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

$$\epsilon = 10^{-7} \quad 10^{-5}$$

Gradient Checking Implementation Notes  $10^{-3}$  - worry

- Don't use in training - only to debug
- If algorithm fails grad check, look at components to try to identify bug.
- Remember regularization  $J(\theta) = \frac{1}{m} \sum L(\hat{y}^{(j)}, y^{(j)}) + \frac{\lambda}{2m} \sum \frac{1}{\ell} \|w^{(\ell)}\|_F^2$
- Doesn't work with dropout  $d\theta$  is gradient of with respect to  $\theta$
- Run at random initialization; perhaps again after some training.

Mini-Batch Gradient Descent,

## Batch vs mini-batch gradient descent

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

baby training sets = mini-batches

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

$(n_{x,m})$        $x^{\{1\}} \ (n_{x,1000})$        $x^{\{2\}} \ (n_{x,1000})$        $\dots$        $x^{\{5000\}} \ (n_{x,1000})$

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

$y^{\{1\}} \ (1,1000)$        $y^{\{2\}} \ (1,1000)$        $\dots$        $y^{\{5000\}} \ (1,1000)$

what if  $m = 5000000$ ?

Mini-batch  $t$ :  $x^{\{t\}}, y^{\{t\}}$

for  $t = 1, \dots, 5000$

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

1 step of gradient descent  
using  $x^{\{t\}}, y^{\{t\}}$   
(as if  $m=1000$ )

Compute cost  $J = \frac{1}{1000} \sum_{i=1}^l (y^{(i)}, y^{(i)})$  from  $x^{\{t\}}, y^{\{t\}}$  +  $\frac{\lambda}{2 \cdot 1000} \sum_l \|w^{(l)}\|_F^2$

Backprop to compute gradients wrt  $J^{\{t\}}$  (using  $(x^{\{t\}}, y^{\{t\}})$ )

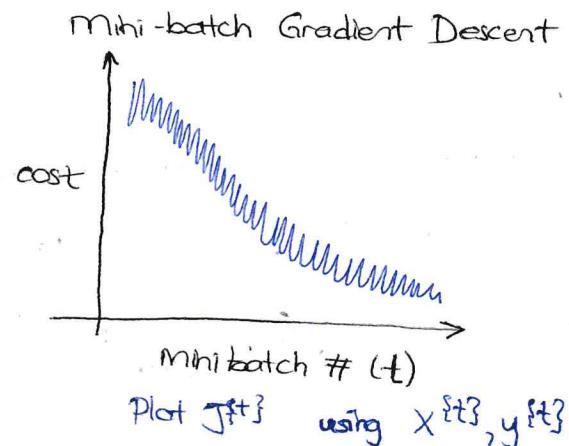
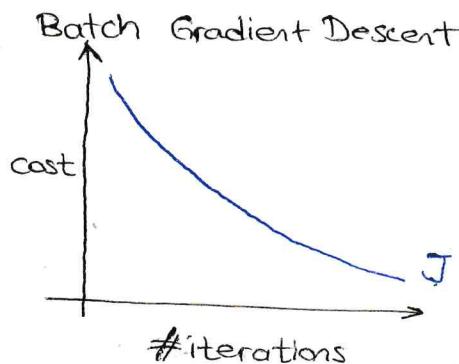
$$w^{(l)} := w^{(l)} - \alpha d w^{(l)}, b^{(l)} := b^{(l)} - \alpha d b^{(l)}$$

"1 epoch"  $\Rightarrow$  1 single pass through 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, one epoch, allows you to take 5000 gradient descent steps.

Understanding Mini-batch Gradient Descent

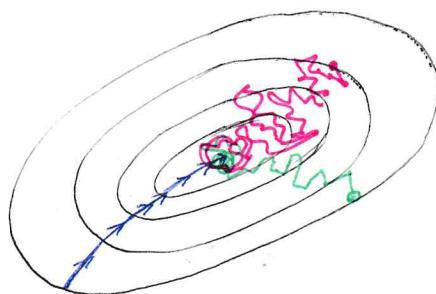
## Training With Mini Batch Gradient Descent



## Choosing your Mini-batch Size

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

If mini-batch size = 1 : Stochastic Gradient Descent. Every example is its own  $(X^{(1)}, y^{(1)}) = (X^{(1)}, y^{(1)}) \dots (X^{(2)}, y^{(2)})$



In practise : Somewhere between 1 and  $m$

## Disadv. of stochastic gradient descent



lose almost all your speed up from vectorization

The way you process each example is going to be very inefficient

In-between  
(mini-batch size not too big/small)



## Fastest learning

- vectorization (a lot of)

- make progress without needing to wait till you process the entire training set.

Disadv. of batch gradient descent  
(mini-batch size =  $m$ )

Too long per iteration

If small training set : Use batch gradient ( $m \leq 2000$ )

Typical mini-batch sizes :

$\underbrace{64}_{2^6} \rightarrow \underbrace{128}_{2^7}, \underbrace{256}_{2^8}, \underbrace{512}_{2^9}, \underbrace{1024}_{2^{10}}$   
more common

Make sure mini-batch fit in CPU/GPU memory

$X^{(t)}, y^{(t)}$

Exponentially Weighted Averages

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

$\alpha \beta = 0.9$  :  $\approx 10$  days' temperature

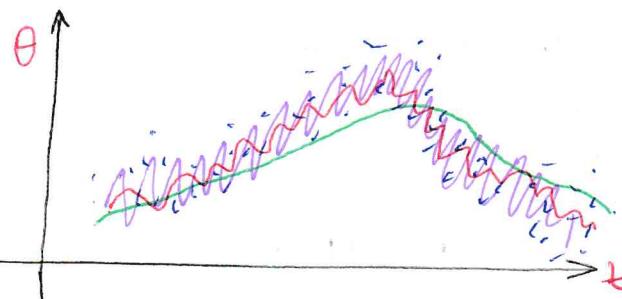
$\alpha \beta = 0.98$  :  $\approx 50$  days

$\alpha \beta = 0.5$  :  $\approx 2$  days

$V_t$  as approximately averaging over

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

temperature

Understanding Exponentially Weighted Averages

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

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

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

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

...

$$\begin{aligned} V_{100} &= 0.1 \cdot \theta_{100} + 0.9 \cancel{V_{99}} (0.1 \cdot \theta_{99} + 0.9 V_{98}) \\ &= 0.1 \theta_{100} + 0.1 \cdot 0.9 \cdot \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} + 0.1 (0.9)^4 \cdot \theta_{96} + \dots \end{aligned}$$

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

$$\frac{(1-e)^{10}}{0.9} = \frac{1}{e}$$

$$\approx \frac{1}{1-\beta}$$

$$\epsilon = 1 - \beta$$

$$\epsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}$$

$$V_0 := 0$$

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

$$V_0 := \beta V_0 + (1-\beta) \theta_2$$

⋮

$$V_0 = 0$$

Repeat {

Get next  $\theta_t$

$$V_0 := \beta V_0 + (1-\beta) \theta_t$$

}

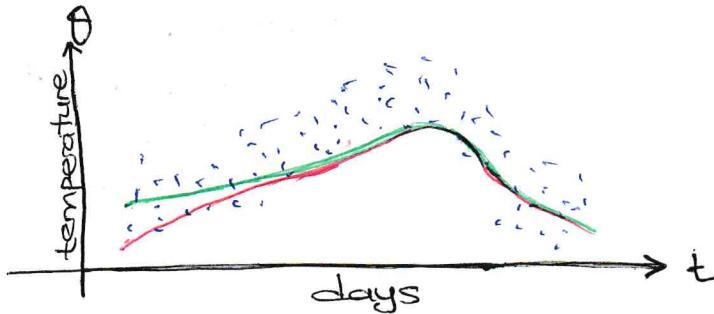
Advantage:- It takes very little memory  
- One line of code

If you were to compute a moving ~~window~~ window, where you explicitly sum over the last 10 days, last 50 days temperature and just divide by 10 or divide by 50 that gives you better estimate. But disadvantage of that of explicitly keeping all the temperatures around and sum of the last 10 days requires more memory. More complicated to implement and more expensive.

## Bias Correction in Exponentially Weighted Average

(4)

Bias correction makes your computation of these averages more accurately,



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

$$V_0 = 0$$

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

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

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

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

### BIAIS CORRECTION

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

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

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

It helps you go from red line to green line

But we are concerned about the bias during this initial phase.

Then bias correction can help you get a better estimate early on.

## Gradient Descent With Momentum



Implement gradient descent with momentum

Momentum:

On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) dW \quad "V_\theta = \beta V_\theta + (1-\beta) \theta_t"$$

$$V_{db} = \beta V_{db} + (1-\beta) db \quad \begin{matrix} \text{acceleration} \\ \uparrow \text{velocity} \end{matrix}$$

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

So what this does is smooth out the steps of gradient descent

### Implementation details

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

On iteration  $t$ :

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

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

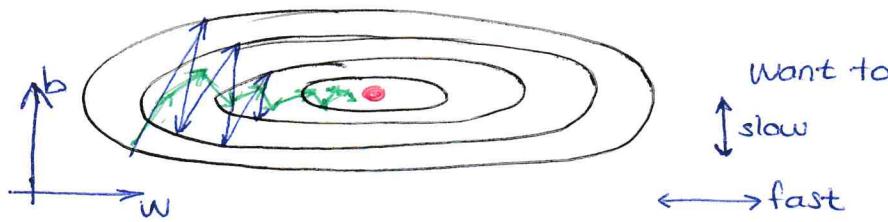
$$V_{db} = \beta V_{db} + (1-\beta) db$$

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

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

$$\frac{V_{dw}}{1-\beta^t}$$

In practise, I don't really see people bothering with bias correction when implementing gradient descent or momentum

RMS prop (Root Mean Square Prop)

On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

$$Sdw = \beta Sdw + (1-\beta) dW^2 \quad \begin{matrix} \leftarrow \text{element} \\ \text{wise} \end{matrix}$$

$$Sdb = \beta Sdb + (1-\beta) db^2 \quad \begin{matrix} \leftarrow \text{small} \\ \rightarrow \text{large} \end{matrix}$$

$$W := W - \alpha \frac{dW}{\sqrt{Sdw + \epsilon}} \quad b := b - \alpha \frac{db}{\sqrt{Sdb + \epsilon}} \quad \begin{matrix} \rightarrow \text{if } Sdw \text{ or } Sdb \text{ is } 0 \\ \epsilon = 10^{-8} \end{matrix}$$

You can use larger learning rate ( $\alpha$ ) and get faster learning without diverging in the vertical direction.

Adam Optimization Algorithm

Adam optimization algorithm is basically taking momentum and rms prop and putting them together.

$$Vdw = 0, Sdw = 0 \quad Vdb = 0, Sdb = 0$$

On iteration  $t$ :

Compute  $dW, db$  using current mini-batch

$$Vdw = \beta_1 Vdw + (1-\beta_1) dW, Vdb = \beta_1 Vdb + (1-\beta_1) db \quad \begin{matrix} \leftarrow \text{"momentum"} \\ \beta_1 \end{matrix}$$

$$Sdw = \beta_2 Sdw + (1-\beta_2) dW^2, Sdb = \beta_2 Sdb + (1-\beta_2) db^2 \quad \begin{matrix} \leftarrow \text{"RMS prop"} \\ \beta_2 \end{matrix}$$

$$Vdw_{\text{corrected}} = Vdw / (1 - \beta_1^t), Vdb_{\text{corrected}} = Vdb / (1 - \beta_1^t)$$

$$Sdw_{\text{corrected}} = Sdw / (1 - \beta_2^t), Sdb_{\text{corrected}} = Sdb / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{Vdw_{\text{corrected}}}{\sqrt{Sdw_{\text{corrected}} + \epsilon}}, b := b - \alpha \frac{Vdb_{\text{corrected}}}{\sqrt{Sdb_{\text{corrected}} + \epsilon}}$$

Hyperparameters

$\alpha$ : need to be tuned

$\beta_1$ : 0.9  $\rightarrow$   $(dW)$   $\rightarrow$  is computing mean of the derivatives

$\beta_2$ : 0.999  $\rightarrow$   $(dW^2)$   $\rightarrow$  is used to compute exponentially weighted average of  $dW^2$

$\epsilon$ :  $10^{-8}$  (doesn't need)

In the typical implementation of Adam, you do implement bias correction.

This is commonly used learning algorithm that is proven to be very effective for many different neural networks of very wide variety of architectures.

Learning Rate DecaySlowly reduce  $\alpha$ 

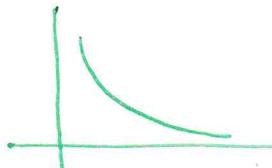
(6)

1 epoch = 1 pass through data

$$\alpha' = \frac{1}{1 + \text{decay-rate} \cdot \text{epoch-num}} \cdot \alpha_0$$

↙  
another hyperparameter

Epoch	$\alpha$	$\alpha_0 = 0.2$
1	0.1	
2	0.067	decay-rate = 1
3	0.05	
4	0.04	
⋮	⋮	

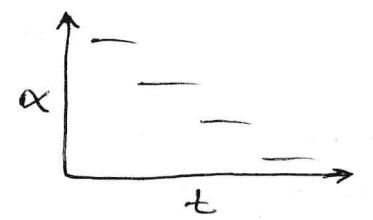

Other Learning Rate Decay Methods

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

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

↑  
mini-batch number

manual decay → this works only if you're training only a small number of models



discrete staircase

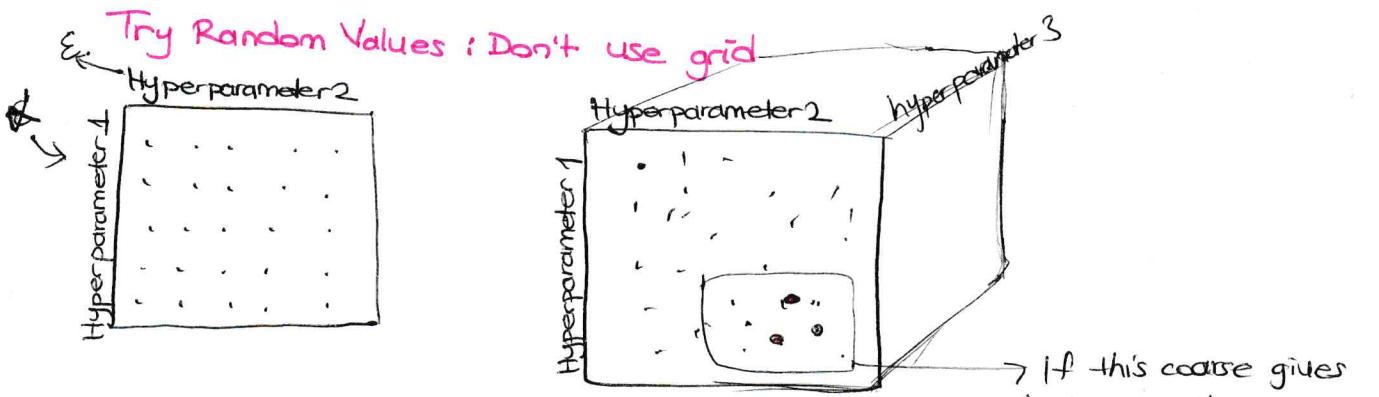
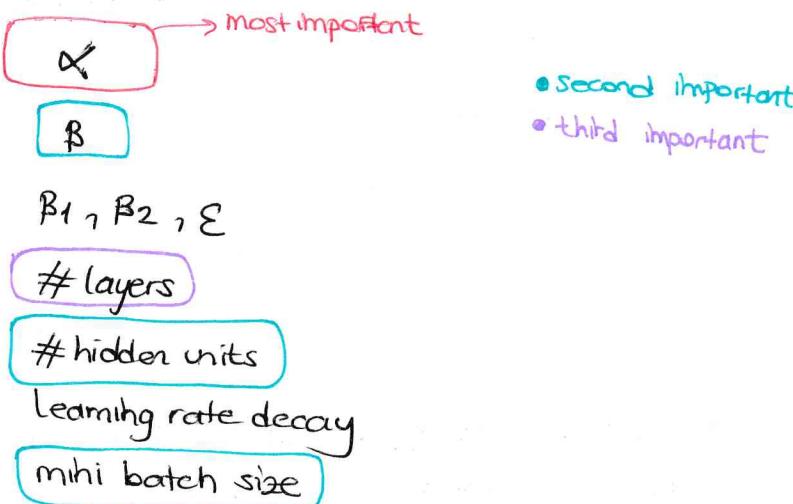
The Problem Of Local Optima

derivative = 0 that point is called a saddle point

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

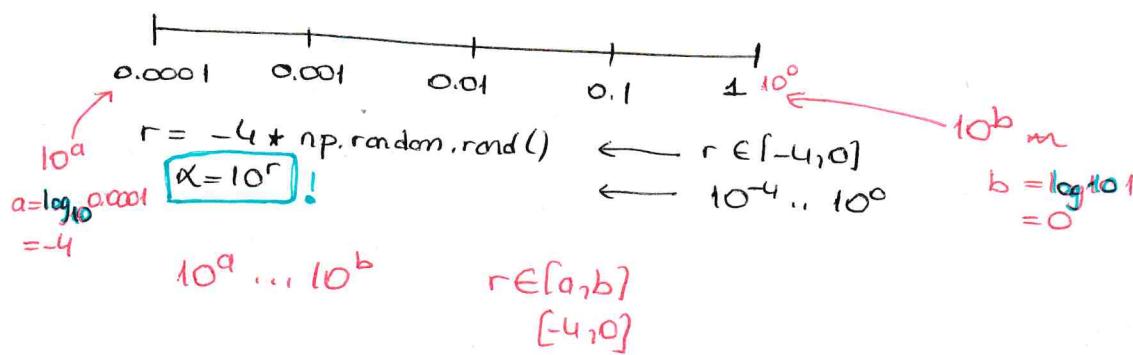
# HYPERPARAMETER TUNING

## Tuning Process



- use random sampling
- adequate search
- (optionally) consider implementing a coarse to fine search process.

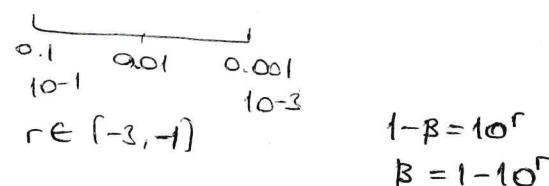
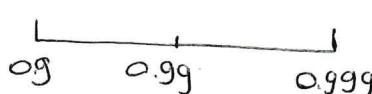
## Using an Appropriate Scale to Pick Hyperparameters

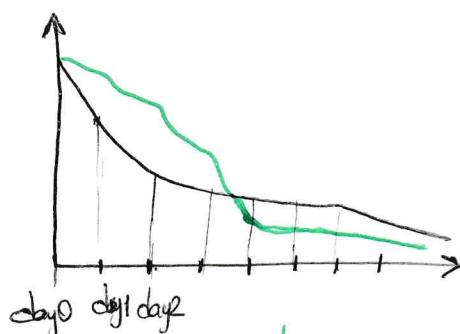
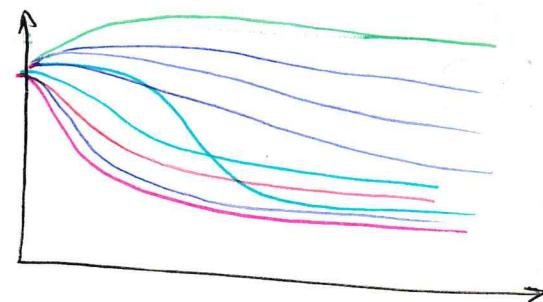


0.9                    0.999

→ It doesn't make sense to sample on the linear scale

$$1-\beta = 0.1 \dots 0.001$$



Hyperparameters Tuning in Practice: Pandas vs CaviarBabysitting One ModelTraining Many Models in Parallel

The way to choose between these two approaches is a function of how much computational resources you have.

For example:

Some online advertising settings as well as in some computer vision applications where there's just so much data, models are so big that it's difficult to train a lot of models at the same time.

If you have enough computers to train a lot of models in parallel

Normalizing Activations in a Network

Batch normalization makes your hyperparameter search problem much easier, makes your neural network much more robust.

Also will enable you to much more easily train even very deep networks.

Created by two researchers, Sergey Ioffe and Christian Szegedy.

$$\begin{aligned} x_1 &\xrightarrow{w, b} \\ x_2 \\ x_3 \end{aligned}$$

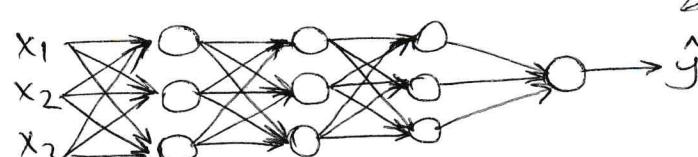
$\mu = \frac{1}{m} \sum_i x^{(i)}$

$x = X - \mu$

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

$x = x / \sigma^2$

↳ normalizing  $x_1, x_2, x_3$   
maybe helps you train w and b more efficiently

Implementing Batch Norm:

Given some intermediate values in NN

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

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

$$\textcircled{3} \quad z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sigma^2 + \epsilon}$$

Wouldn't it be nice if you can normalize the mean and variance of  $a^{(2)}$  to make the training of  $w_3, b_3$  more efficient?

Can we normalize  $a^{(2)}$  so as to train  $w^{(3)}, b^{(3)}$  faster

$$z^{(1)}, \dots, z^{(m)} \xrightarrow{\text{learnable parameters}} z^{(1)(i)}$$

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

↳ learnable parameters of model

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

$B = M$  is true

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

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

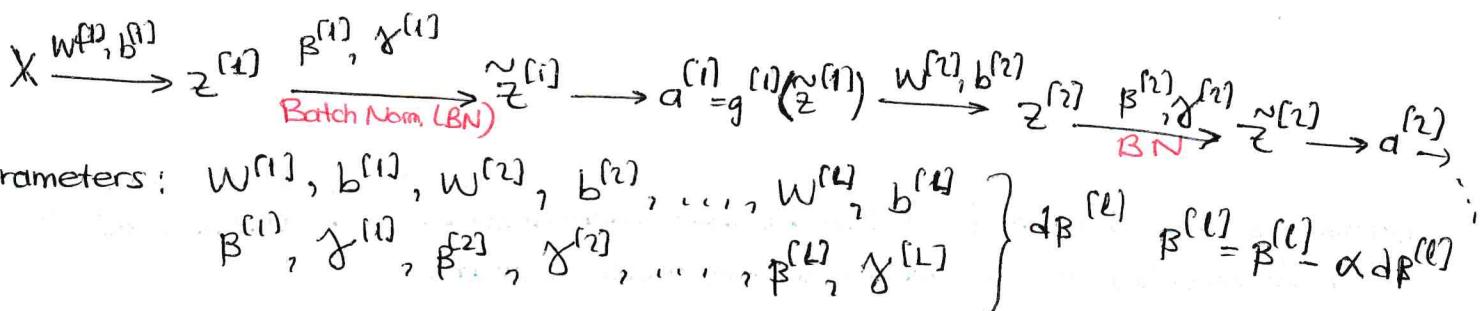
With the parameters gamma and  $\beta$ , you can make that your  $z^{(i)}$  values have the range of values that you want.

What it really does is it normalizes the mean and variance of these hidden units values to have some fixed mean and variance.

could be 0 and 1 or other value

It's controlled by  $\gamma$  and  $\beta$

### Fitting Batch Norm into a Neural Network



this  $\beta$  is different than the  $\beta$  in Adam, RMSprop or momentum.

tf.nn\_batch\_normalization

In practice, Batch norm is usually applied with mini batches of your training set.

Adding any constant to all of the examples in the minibatch, it doesn't change anything. Because any constant you add will get cancelled out by the mean subtraction step. So if you're using batch norm, you can eliminate  $b^{(l)}$ .

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$$

$$\tilde{z}^{(l)} = W^{(l)} a^{(l-1)}$$

$$z^{(l)} = \gamma^{(l)} \tilde{z}^{(l)} + \beta^{(l)}$$

is a parameter that controls that ends up affecting the shift or the biased terms.

Parameters:  $W^{(l)}, b^{(l)}, \beta^{(l)}, \gamma^{(l)}$   
 $(n^{(l)}, 1) \quad (n^{(l)}, 1) \quad (n^{(l)}, 1) \quad (n^{(l)}, 1)$

are used to scale the mean and variance of each row of the hidden units to whatever the network wants to set them to.

Implementing Gradient Descent:

for  $t=1 \dots \text{numMiniBatches}$

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

In each layer, use BN to replace  $z^{(l)}$  with  $\tilde{z}^{(l)}$

Use backprop to compute  $dW^{(l)}, db^{(l)}, d\beta^{(l)}, d\gamma^{(l)}$

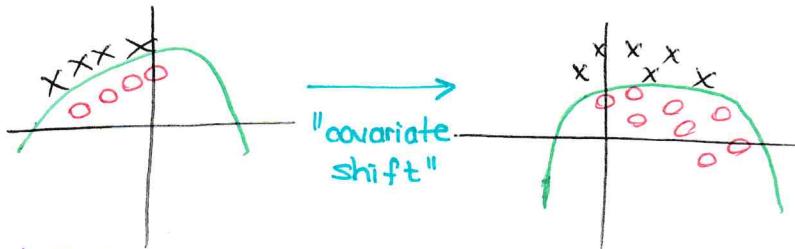
Update parameter  $W^{(l)} := W^{(l)} - \alpha dW^{(l)}$   
 $\beta^{(l)} := \beta^{(l)} - \alpha d\beta^{(l)}$   
 $\gamma^{(l)} := \dots$

Works for all optimization algorithms

Why Does Batch Norm Work?

One intuition behind why batch norm works is, this is doing the similar thing for the values in your hidden units and not just for your input layer.

## X Learning on Shifting Input Distribution:



## Detailed Exp:

Batch norm limits the amount to which updating the parameters in the earlier layers can affect the distribution of values.

So, batch norm reduces the problem of the input values changing, it really causes these values to become more stable, so that the later layers of the neural network has more firm ground to stand on. And so it allows each layer of network to learn by itself, a little bit more independently of other layers and this is the effect of speeding up of learning in the whole network.

The earlier layers don't get to shift around as much, because they're constrained to have the same mean and variance and so this makes the job of learning on later layers easier.

## X 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. (<sup>he thinks</sup> unintended effect) [w128125b](#)
  - $\mu$  and  $\sigma^2$  are noisy
  - small affect (not huge)
  - so you might choose to use batch norm together with dropout

Using with dropout makes more powerful regularization affect

By using <sup>LARGE</sup> minibatch size, you're reducing this noise and therefore also reducing this regularization affect.

mini batch :  $64 \rightarrow 512$

Don't turn to batch norm as a regularization.

Use it as ~~an~~ way to normalize your hidden units activations and therefore speed up learning.

Batch Norm At Test Time

Batch norm processes your data one mini batch at a time, but the test time you might need to process the examples one at a time.

## TRAIN

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

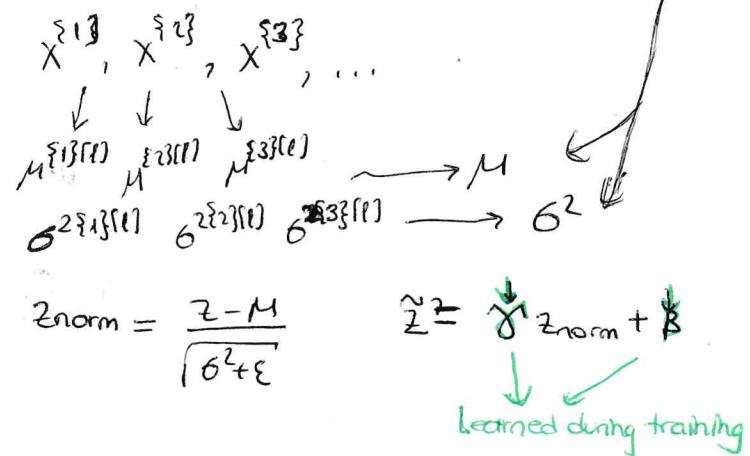
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - M)^2$$

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

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

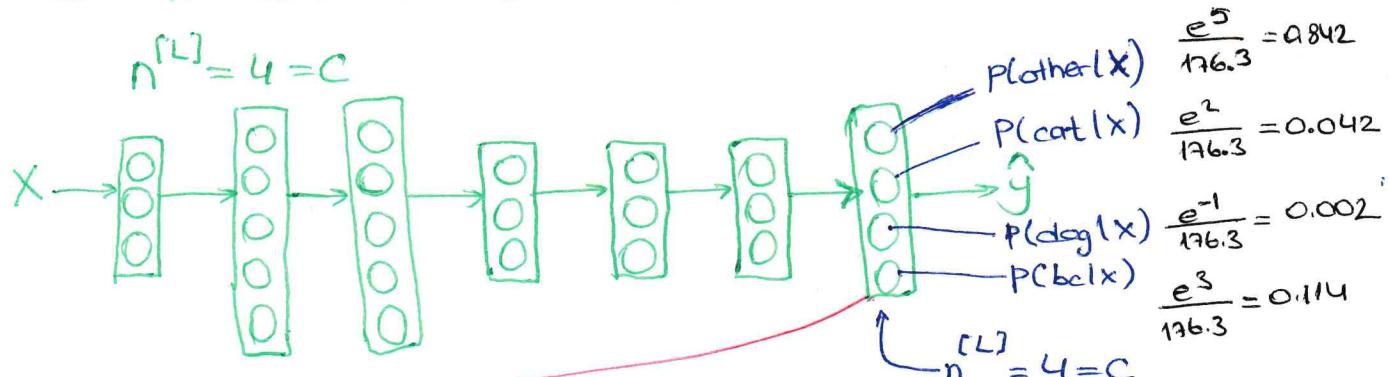
## TEST

$M, \sigma^2$ ; estimate using exponentially weighted average (across minibatches)

Softmax Regression

There's a generalization of logistic regression called Softmax Regression,

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



Softmax Layer:

$$z^{[L]} = W^{[L]} \cdot a^{[L-1]} + b^{[L]} \quad \rightarrow \text{dimension } (4, 1)$$

Activation Function:

$t = e^{z^{[L]}}$	$\dots$	$(4, 1)$
temporary variable	$a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}$	$, a_i^{[L]} = \frac{t_i}{\sum_{j=1}^4 t_j}$

$$a^{[L]} = g^{[L]}(z^{[L]}) \quad (4, 1)$$

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

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

$$\sum_{j=1}^4 t_j = 176.3$$

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

The unusual thing about the Softmax activation function is it takes an input  $(4, 1)$  vector, and outputs a  $(4, 1)$  vector.

Training a Softmax ClassifierUnderstanding Softmax

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

"softmax"

$$a^{[L]} = g^{(L)}(z^{[L]}) = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

"hard max!"

Softmax regression generalizes logistic regression to C classes.

If  $C=2$ , softmax reduces to logistic regression

Loss function

$$(4,1) y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \text{cat} \quad a^{[L]} = \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$y_2 = 1$   
 $y_1 = y_3 = y_4 = 0$

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

$\text{C=4} \quad -y_2 \log y_2 = -\log \hat{y}_2 \quad (\text{make } \hat{y}_2 \text{ big})$

More generally, what this loss function does is it looks at whatever is the ground true class in your training set, it tries to make the corresponding probability of that class as high as possible.

Gradient Descent with Softmax

Backprop :  $\frac{dz^{[L]}}{dJ} = \hat{y} - y$

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

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

$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix} \quad = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \dots$$

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

```
import numpy as np
import tensorflow as tf
w = tf.Variable(0, dtype=tf.float32)
cost = w**2 - 10*w + 25 # tf.add(tf.add(w**2, tf.multiply(-10., w)), 25)
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
init = tf.global_variables_initializer()
session = tf.Session() session.run(init) session.run(w)
```

Backprop is already built in add, multiply and square function in Tensorflow.

```
w = tf.Variable(0, dtype=tf.float32)
```

~~coefficients~~

```
cost = w**2 + 10*w + 25
```

```
train = ...
```

~~init~~

```
session = ...
```

```
session.run(train)
```

```
session = tf.Session()
```

```
session.run(init)
```

```
print(session.run(w))
```

```
coefficients = np.array([[-1.], [-20], [100]])  
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]  
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)  
  
init = tf.global_variables_initializer()  
session = tf.Session()  
train  
session.run(train, feed_dict={x: coefficients})
```

with `tf.Session()` as session:

```
session.run(init)
```

```
print(session.run(w))
```