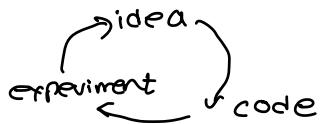


dec1: Train/Dev/Test sets:

ML highly iterative process

- # layers
- # hidden units
- learning rate
- activation function



Data	training set	cross validation dev set	test set
------	--------------	-----------------------------	----------

previous : 70% - train 60 - train $m = 10,000 \text{ to } 100,000$
 30% - test 20 - dev
 20 - test

Big data : $m = 1,000,000$
 98% train 99.5% train
 1% dev 0.4% dev
 1% test 0.1% test

Mismatched train/test distribution:-

Training set:

cat pictures from webpage
 - high res
 - good

Dev/test sets:

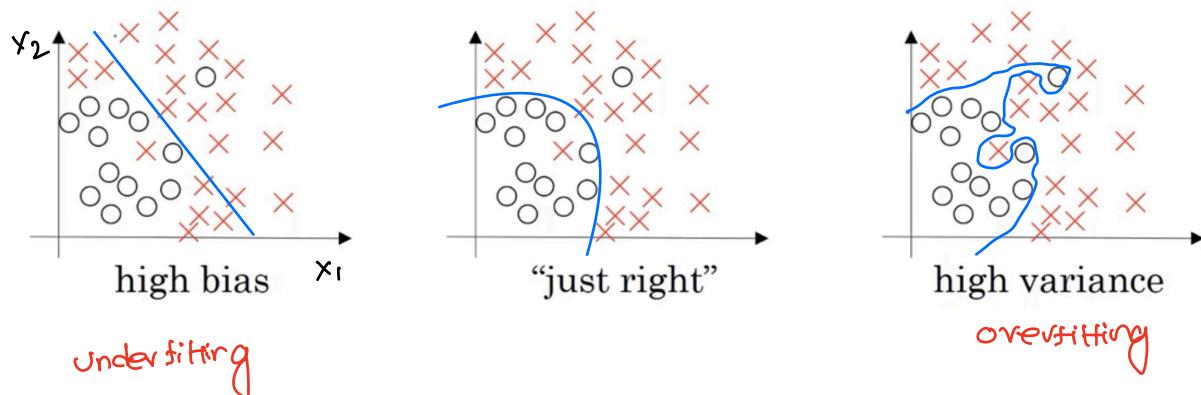
cat pictures from phone app
 - low res
 - blurry

Make sure dev, test come from same distribution

- Not having a test set might be okay. (only train, dev test)

↓
test

dec2: Bias and Variance:



Cat classification



Human error $\approx 0\%$.
Optimal Bayes error $\approx 0\%$.

Train set error : 1 %.

Dev set error : 11 %.

high variance

15 %

16 %
high bias

15 %

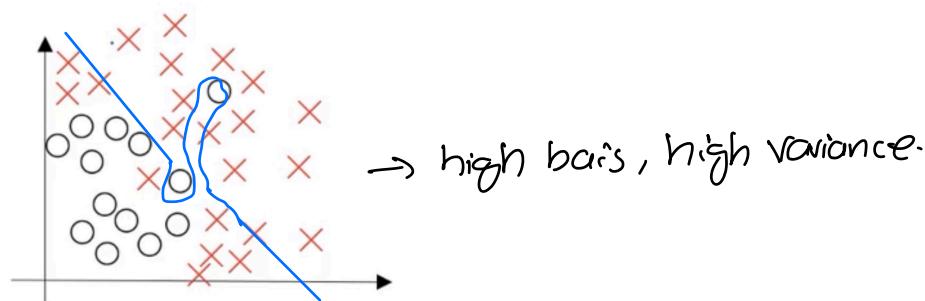
30 %
high bias
high variance

0.5 %

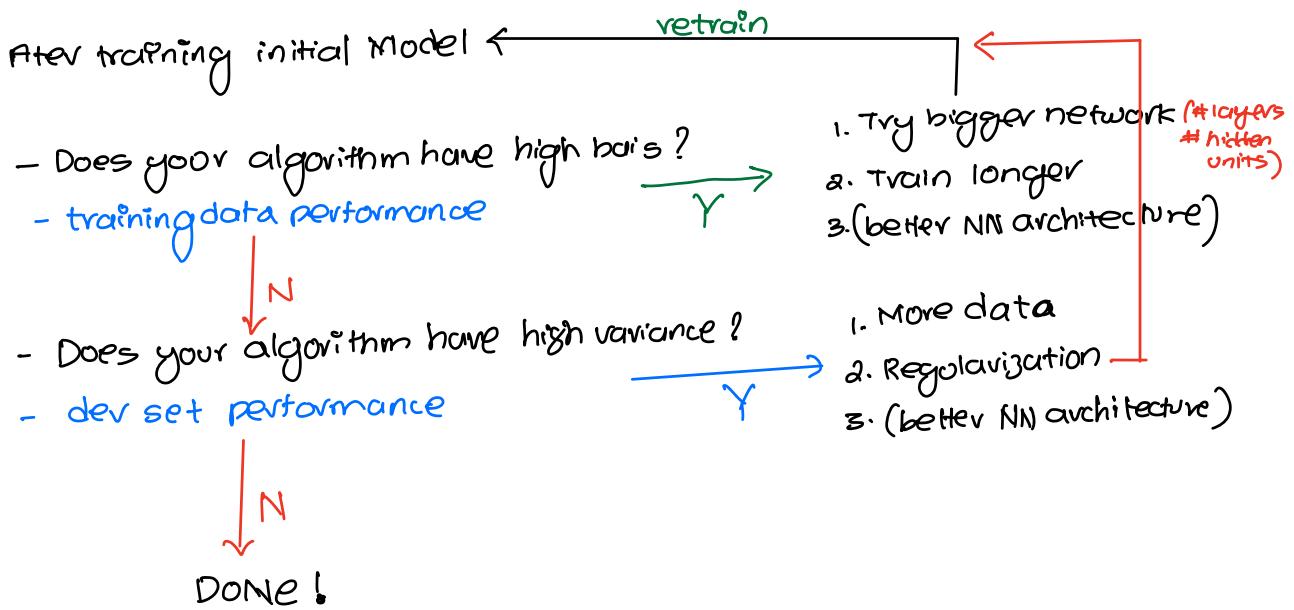
1 %
low bias
low variance

if optimal / Bayes error is 15% , such as blurry images .

then may train 15% , dev 16% - low bias, low variance .



dec3. Basic Recipe for Machine learning:



Bias Variance Tradeoff:

Tradeoffs: previous era: ↑ bias ↓ variance bias ↑ variance

Big data era: Move data, Bigger Network → WORKS
↓
Deep learning reduces bias independently
reduces variance

Lec4: Regularization:

- Use when model is overfitting or high variance problem.
- use get more training data

Logistic Regression:-

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

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

$$\alpha_2 \text{ regularization} = \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$\alpha_1 \text{ regularization} = \frac{\lambda}{2m} \|w\|_1 = \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j|$$

λ = regularization parameters = lambda

Neural Network Reg:

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

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

- Frobenius Norm.

Before

$$d\omega^{(l)} = \text{from backprop}$$

$$\omega^{(l)} = \omega^{(l)} - \alpha d\omega^{(l)}$$

After reg

$$d\omega^{(l)} = (\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)}$$

$$\omega^{(l)} = \omega^{(l)} - \alpha d\omega^{(l)}$$

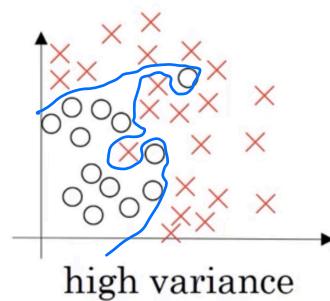
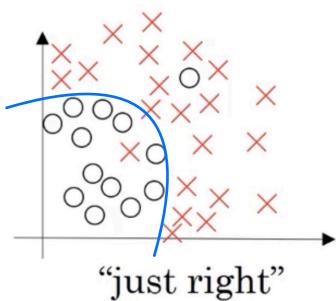
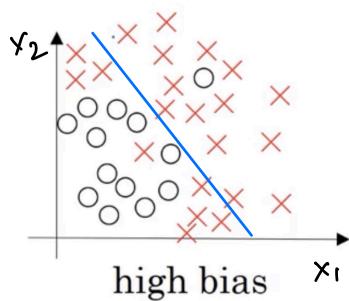
weight decay $\omega^{(l)} = \omega^{(l)} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)} \right]$

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

$$= \left[1 - \alpha \frac{\lambda}{m} \right] \omega^{(l)} - \alpha d\omega^{(l)}$$

weight decay. < 1

Dec 5: Why does Regularization reduces Overfitting:-

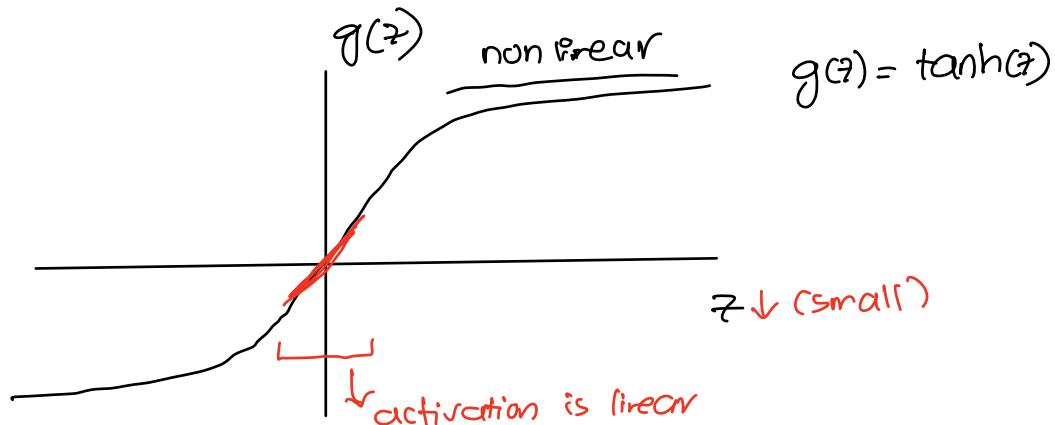


Cost with Reg:-

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

if $\lambda = \infty$ then to min J $\omega^{(l)} \approx 0$

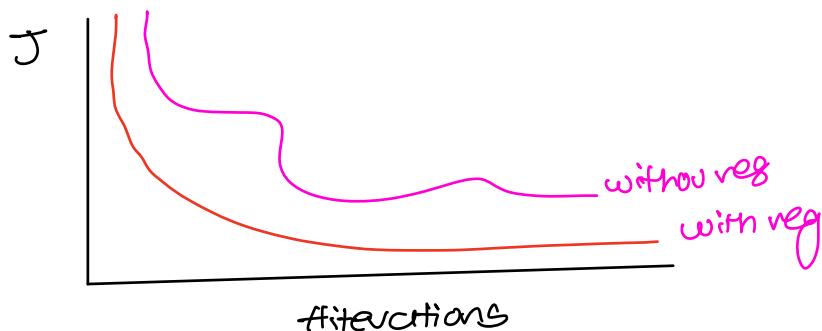
\Rightarrow Neural Network becomes much simpler network.



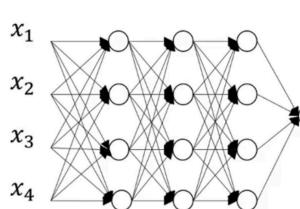
$$\lambda \uparrow \quad w^{(l)} \downarrow \quad z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)} \Rightarrow z^{(l)} \downarrow$$

$\downarrow g(z^{(l)}) \Rightarrow$ roughly linear regression.
 \Rightarrow deep network \Rightarrow linear regression
 \Rightarrow model become simpler.

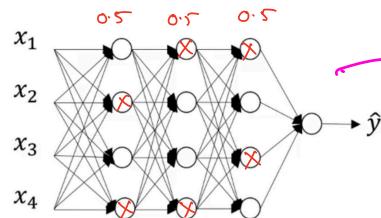
$$J(\dots) = \sum_i \alpha(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{(l)}\|_F^2$$



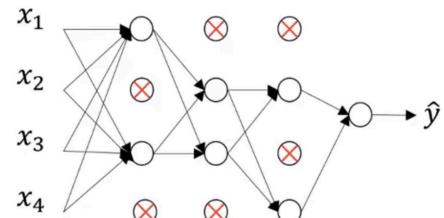
dec 6: Dropout Regularization!



NN-1



NN-2



smaller Network

if NN-1 overfits your data, then we can change following

- do through each of layers in NN-2
- set some probability of eliminating a node in NN-2

\Rightarrow smaller Network will generalize data better, and reduce the overfitting.

Implementing dropout ("Inverted dropout") = /keepdrop

Illustrate with layer $l=3$

$\text{keepprob} = 0.8$
 $d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep-prob}$

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

$a_3 /= \text{keep-prob}$ makes test time easier (less scaling problem).

if $a_3.\text{shape} = (50, m) = (n^{[3]}, m)$

with dropout: keep prob = 0.8

\Rightarrow on avg layer 3 has 10 units shut off $= 50 * 0.2$

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

in order not to reduce the Expected value of $z^{[4]}$, we $a^{[3]} /= \text{keep prob}$

Predictions at test time

$$a^{(0)} = x$$

NO dropout at test time } No random predictions.

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$$

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

$$z^{(2)} = w^{(2)} a^{(1)} + b^{(2)}$$

$$a^{(2)} = g^{(2)}(z^{(2)})$$

Q7: Why does dropout work?

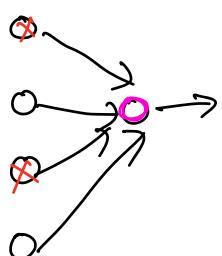
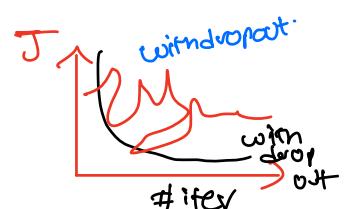
Intuition: can't rely on any one feature, so have to spread out weights

dropout randomly knock out hidden units

every iteration working with smaller network

have regularizing effect

} previous intuition.



- can't rely on any one feature
→ hence spreads weights out
→ this will shrink weights → ℓ_2 regularization.
→ prevents overfitting

dropout formally can be shown to be adaptive form
 ℓ_2 regularization.

$$\omega^{(1)} = (7, 3) \quad \boxed{\omega^{(2)} = (7, 7)}, \quad \omega^{(1)} = (3, 7), \quad \omega^{(2)} = (2, 3) \quad \omega^{(1)} = (1, 2)$$

Keepprob = 0.5

$0.7 \quad 0.7 \quad 1.0 \quad 1.0$

↓

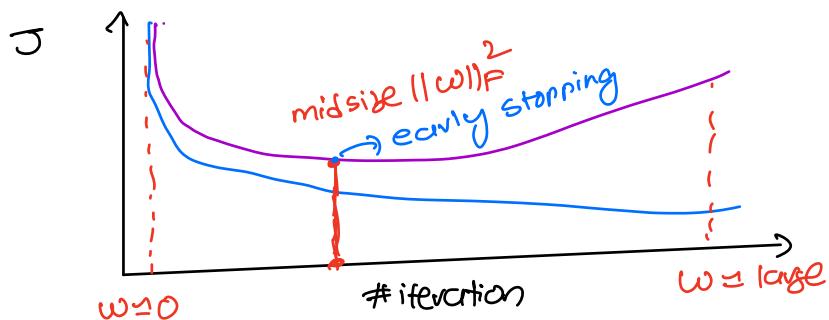
less with more parameters

decs: Other Regularization:-

(i) Data Augmentation:-

- (i) flip horizontally images (double training set)
- (ii) random crops
- (iii) random rotation, distortion

(ii) Early stopping:- (overfitting)



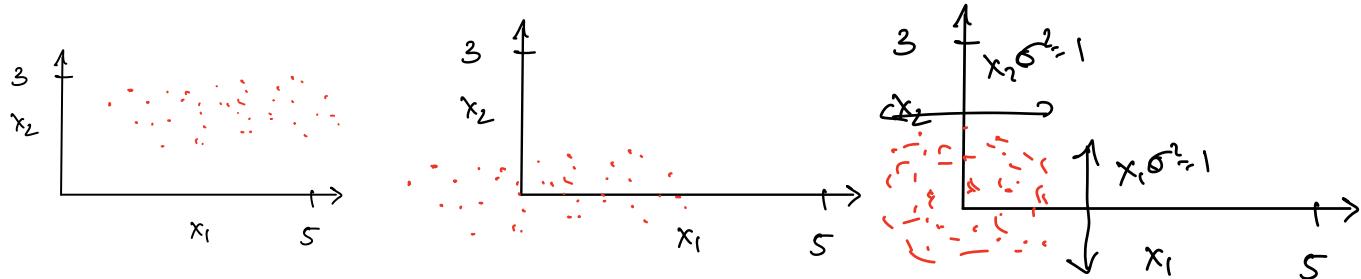
problem? (orthogonalization is not possible)

- ① optimize cost function J
- gradient descent
 - momentum
- ② Not overfit
- regularization
 - more data
- $\min_{w,b} J(w,b)$
- separate task.
orthogonalization.
- reduce variance

Setting UP Your Optimization Problem:

deca: Normalizing inputs:-

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



① subtract mean

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

$$x := x - \mu$$

② Normalize variance

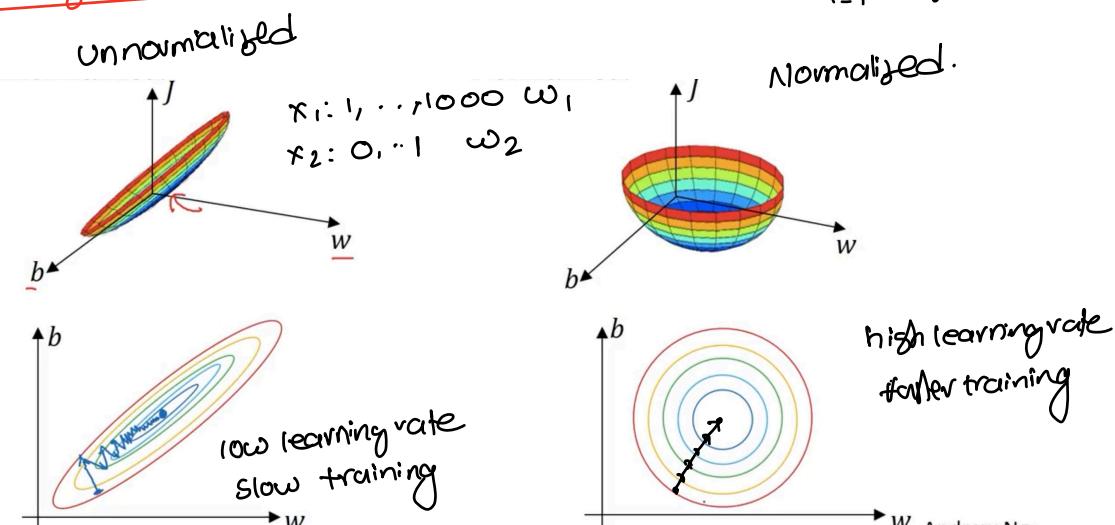
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i) \times 2}$$

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

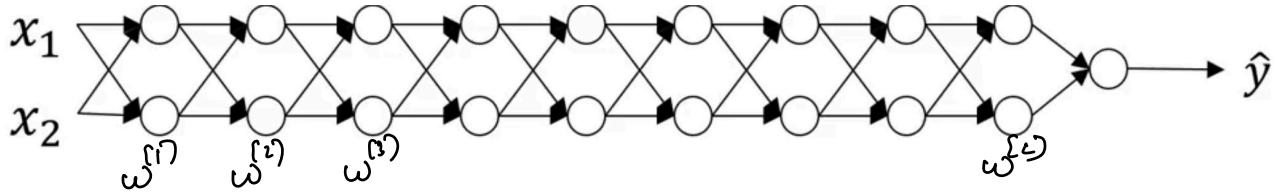
use same μ, σ for test set same as training set

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

Why Normalize inputs?



lec10: Vanishing /Exploding gradients: (DeepNN problem)



$$g(z) = z \quad (\text{linear activation})$$

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

$\boxed{z = \sum_j w^{[j]} a^{[j]}}$

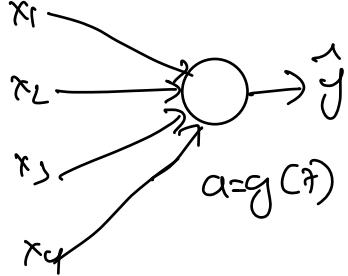
$z = \sum_j w^{[j]} a^{[j]}$
 $a^{[j]} = g(z^{[j]}) = z$

$$L=1 \quad \omega^{[1]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \quad \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \quad \approx \text{Very large}$$

$$L=1 \quad \omega^{[1]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \quad \hat{y} = w^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x \quad \approx \text{Very low}$$

\rightarrow same logic d ω , applies

deciII. Weight Initialization for Deep Networks:-



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b^{\theta}$$

large $n \rightarrow$ smaller w_i

$$\text{var}(w) = \frac{1}{n} \quad \text{if } g(z) = \text{relu}$$

$$\text{var}(w) > \frac{2}{n}$$

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

if $g(z) = \text{relu}$

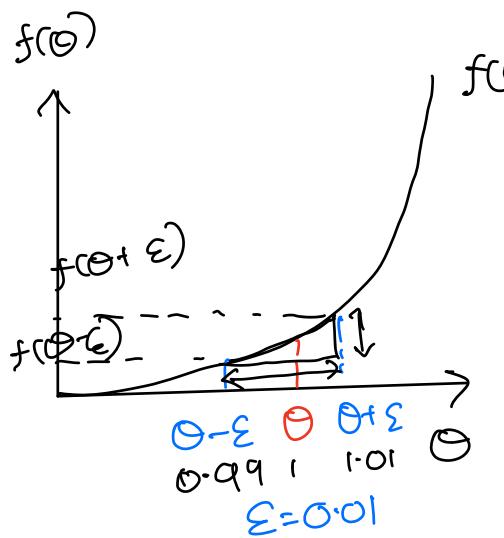
$$w^{(l)} = np.random.rand(\text{shape}) * np.sqrt\left(\frac{d}{n^{(l-1)}}\right)$$

if $g(z) = \tanh$

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

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

Dec 2: Numerical approximation of gradients



$$m = \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{(\theta + \epsilon) - (\theta - \epsilon)}$$

$$m = \frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 30001$$

$$g(\theta) = 3\theta^2 = 3(1) = 3$$

approx error = 0.0001

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

Dec 3: Gradient checking for NN:

take $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}$ and convert to 10^{-3}

big vector θ

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

take $d\omega^{(1)}, db^{(1)}, \dots, d\omega^{(L)}, db^{(L)}$ and reshape into

big vector $d\theta$

is $d\theta$ gradient of $J(\theta)$?

for each i :

$$d\theta_{approx}^{(i)} = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$
$$\simeq d\theta^{(i)} = \frac{\partial J}{\partial \theta^i}$$

$$d\theta_{approx} \simeq d\theta \quad ? \quad \epsilon = 10^{-7}$$

check

$$\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \simeq 10^{-7} - \text{great!}$$
$$\text{if } 10^{-5} - \text{check}$$
$$\text{if } 10^{-3} = \text{bug}$$

dec14: Gradient checking implementation!

1. don't use in training - only to debug
 2. if algorithm fails grad check, look at components to try to identify bug
 3. remember regularization
- $$J(\theta) = \frac{1}{m} \sum_{i=1}^m d(C_j^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|\theta\|_F^2$$
4. doesn't work with dropout. ($\text{keepprob} = 1.0$) then check
 5. run at random initialization; perhaps again after some training