

Course 2: Improving Deep Neural Networks: Hyperparameter Tuning, Regularization and Optimization

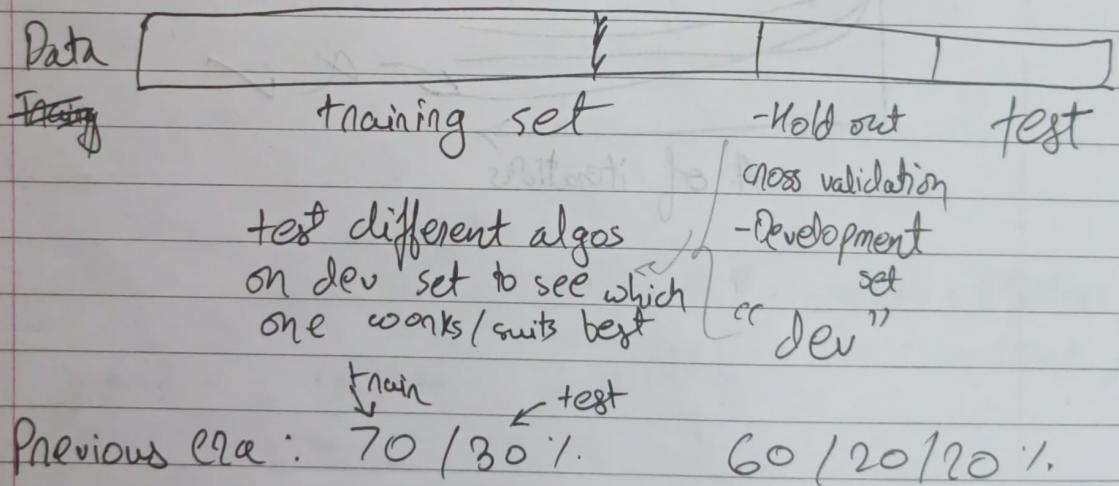
Week 1 : Setting up your ML Application

1) Train / Dev / Test Sets :

Applied ML is a highly iterative process as initially you do not know # layers, # hidden units, learning rate, activation function, etc ...

DL is used in NLP, CV, Structured data

logistics
Ads search Security



Big data era: 1,000,000 10,000 10,000

98, 1, 1 %.

Mismatched train / test distribution

eg: ✓

✓

Training set:

Dev/Test sets:

(cat pictures from web pages) \longleftrightarrow (cat pictures from user using your app)

~~Make sure that dev and Test sets come from the same distribution~~

Note: Not having a test set might be okay
ie. only train and dev set.

2) Bias and Variance

eg: Assuming Human error in classifying cat/not cat $\approx 20\%$.
 Train set error: 1%
 Dev set error: 11% (relatively poor)

\rightarrow Overfitting the training set i.e. high variance.

Train set error: 15%.

Dev set error: 16%.

\rightarrow Underfitting the data (high bias).

Training set error: 15%.

Dev set error: 30%.

\rightarrow Both high variance and high bias.

Training set error: 0.5%.

Dev set error: 1%.

\rightarrow Low variance and Low bias.

We classify bias and variance depending on
Optimal (Bayes) error i.e. human error

But if Bayes error is higher than train set error then algo won't be considered as high bias

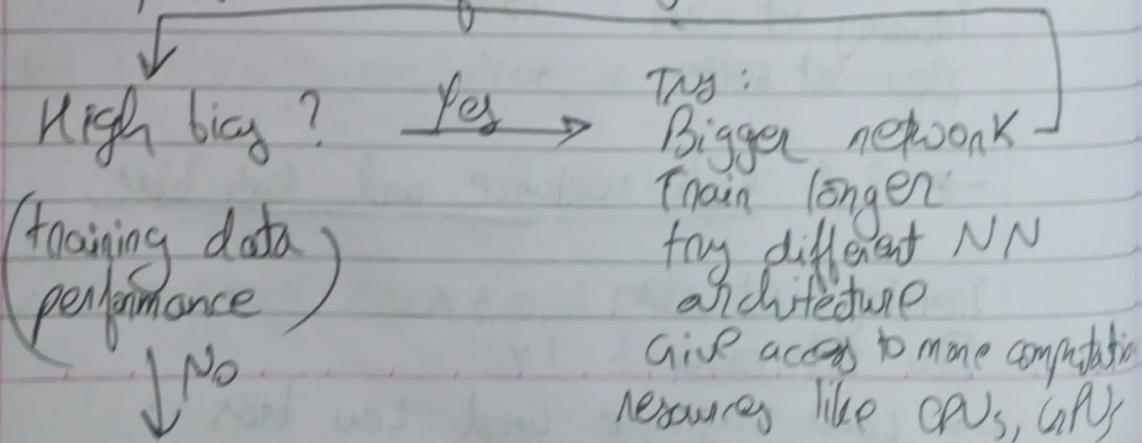
e.g.: Bayes error = 15%. ↗ ~~high~~ low bias
Train set = 15%. ↗

→ when blurry images are shown so then humans as well as machines cannot classify correctly

Note: Linear classifier has high bias
(under fits data)

3) Basic 'Recipe' for ML

→ tells us whether the algo has bias or variance problem on both, hence improves the performance of our algorithm.





High Variance?
(dev set performance)

Yes

→ get More data
but more expensive

Try Regularization
(Reduces overfitting)

No

diff. NN architecture

Done

pre-deep learning : Bias Variance trade off

↑ ↓

↓ ↑

Big data era : They both can be independently reduced without hurting each other using above methods.

Regularization :

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

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

for regularization :

γ regularization parameter

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

omit b just a no.

lambda → reserved keyword in python
 $\therefore \text{we use } \text{use}(\text{lambda}) \rightarrow \text{hyperparameter}$

norm of w squared

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

$$L_1 \text{ regularization : add } \frac{\lambda}{m} \sum_{j=1}^m |w_j| = \frac{\lambda}{m} \|w\|_1,$$

→ w will be sparse
 (ie it will have lot of zeros in)

This will ~~compres~~ our model and hence
~~needs~~ less memory to store the model.
 But this helps only a little bit hence we
 do not use L_1 neg.

Neural Network - $J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]})$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

Adding regularization :

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

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

"Frobenius norm" → sum of squares of elements
 in the matrix

Gradient descent :

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

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

L_2 reg \rightarrow "weight decay"

$$\begin{aligned} w^{[l]} &= w^{[l]} - \alpha \left[(\text{from back}) + \frac{\lambda}{m} w^{[l]} \right] \\ &= w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from back}) \end{aligned}$$

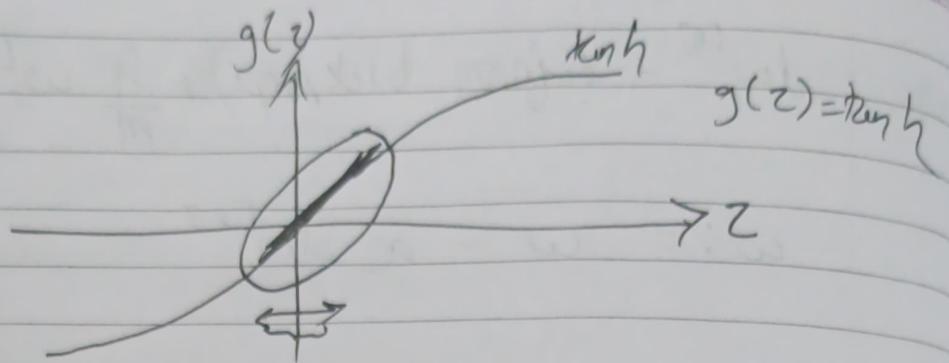
$$\underbrace{\left(1 - \frac{\alpha \lambda}{m}\right) w^{[l]}}_{< 1}$$

This term effectively shrinks the weight w , on each iteration.

The regularization term $\frac{\lambda}{m} w$ increases the

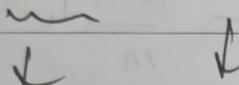
effective cost for large weights (by appropriately selecting a large value of λ , thereby penalizing them).

This will in turn prevent overfitting as the NN will become much simpler as weights are now small hence they won't contribute much hence almost eliminating them and given the NN a structure of logistic regression.



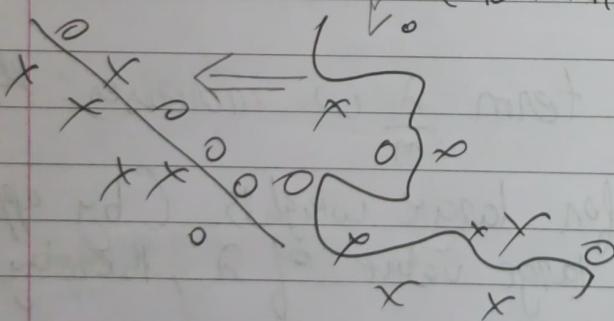
if $\lambda \uparrow$ then $w^{[e]} \downarrow$ as it is penalized due to cost fn

$$\therefore z^{[e]} = w^{[e]} a^{[e-1]} + b^{[e]} \quad \text{grad. descent}$$



∴ Every layer \approx Linear

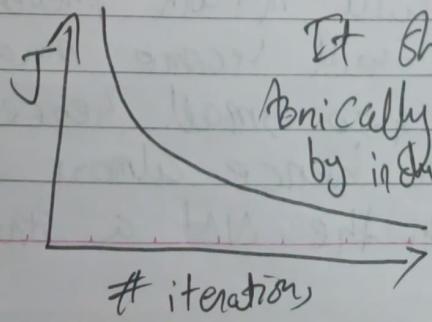
∴ The complex Deep NN just becomes a linear regression, which will only be able to compute a linear function, so it won't be able to fit complicated decisions like (non linear decision boundaries)



Hence overfitting is prevented by making complex non-linear f's to linear f's.

To debug gradient descent:

Plot cost J vs # iterations



It should monotonically decrease by including term $\frac{\lambda}{m} \sum w_j^2$

→ Helps to reduce variance of model

Page No.

Date

Dropout regularization :

Implement through "Inverted dropout" (^{most common})

Illustrate with layer $l = 3$, keep-prob = 0.8

Probability that a unit is not eliminated from l .

$d_3 = np.random.rand(a_3.shape[0], a_3.shape[1]) < 0.8$

$a_3 = np.multiply(a_3, d_3)$

* $d_3 / \text{keep-prob}$ \rightarrow inverted dropout technique

50 units \rightsquigarrow 10 units eliminated ($: 0.8 \times 50 \approx 10$)

$$z^{(4)} = w^{(4)} a^{(3)} + b^{(4)}$$



Reduced by 20% if we do not write *

Hence to compensate for the loss of hidden units and so that the output $z^{(4)}$ does not change we shall divide a_3 by 0.8

i.e. instead the expected value of a_3 and in turn $z^{(4)}$ remains the same as it should be ~~with~~ if all hidden units were present.

< 0.8
implemented in python as:

$$d_3 = (a_3 < \text{keep-prob}).astype(\text{int})$$

Making predictions at test time

$$a^{[0]} = x$$

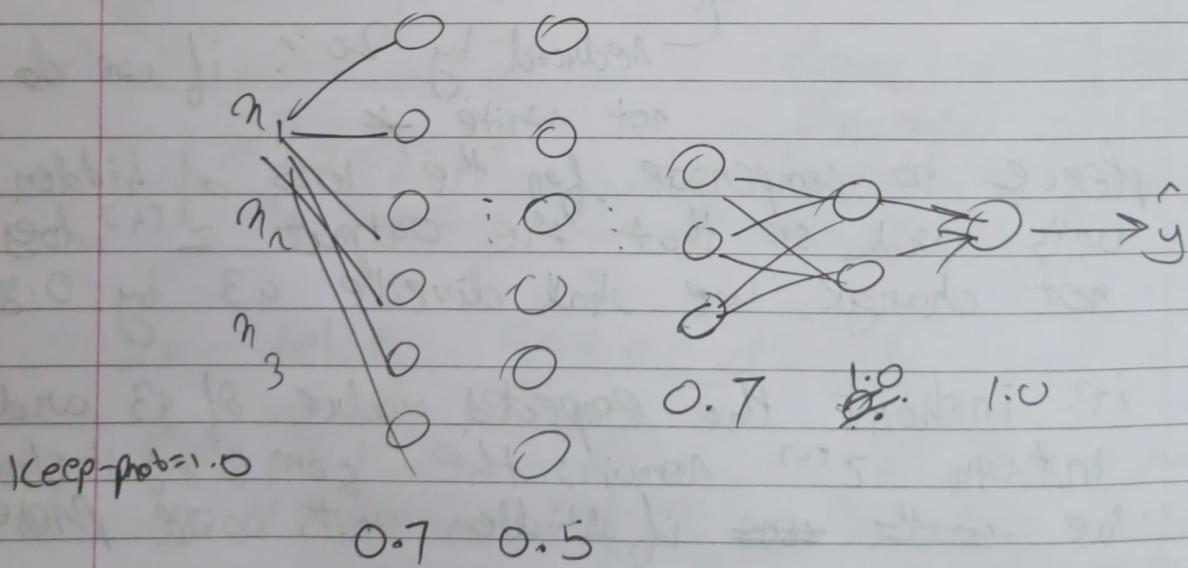
No drop out at test time

$$\begin{aligned} z^{[1]} &= w^{[1]} a^{[0]} + b^{[1]} \\ a^{[1]} &= g^{[1]}(z^{[1]}) \end{aligned}$$

$$\begin{aligned} z^{[2]} &= w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= g^{[2]}(z^{[2]}) \end{aligned}$$

...
↓
g

We shall implement drop out during training and not testing as it would add noise to the predictions if we use it in test time.



Keep the values of keep-prob low for those layers where many hidden units are not required which

complexes the NN overall.

Dropout is used in CV as input images size is big so we do not have enough data to store them.

1) Data Augmentation

Other Regularization tech. to reduce overfitting

Ex : Recognize Cats / Non cats

Getting data from outside can be expensive

so instead we can take an image and make a copy of it by rotating it to get a new image, so this would save us a lot. Also we can zoom into cat images and rotate them by a few degrees to make new combination of images. This also reduces overfitting.

Ex : Number recognition (optical characters)

From the above example, we can implement similar methods by distorting images, rotating images, etc.

2) Early stopping



Stopping midway when dev set error is least hence do not continue the grad. descent further and extract mid-way w^F .

Downside:

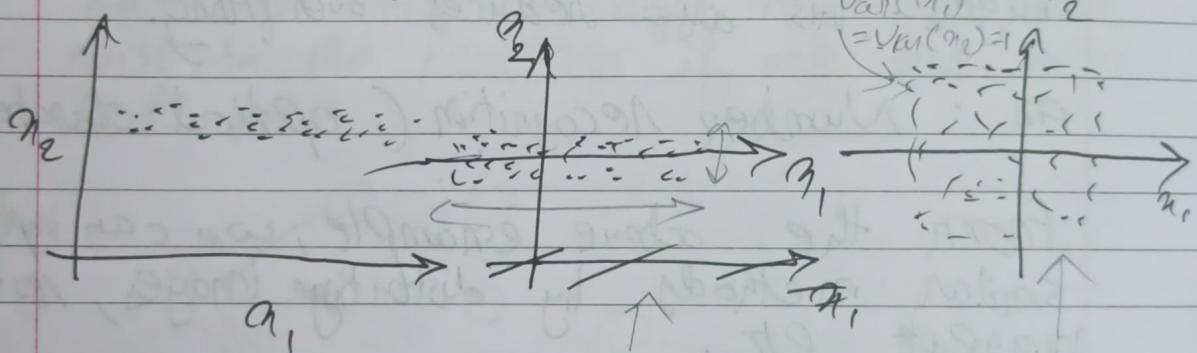
By using early stopping, we are not fully using gradient descent so that is not going to fully minimize the cost function.

Alternatively we can use L₂ Reg. instead of early stopping.

Setting up your optimization problem

① Normalizing Inputs

e.g.: To input features a_1 & a_2 :



Subtract mean:

Normalize variance

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

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

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

element-wise square

$$\lvert n \rvert = \sigma$$

Here $\text{Var}(a_1) > \text{Var}(a_2)$

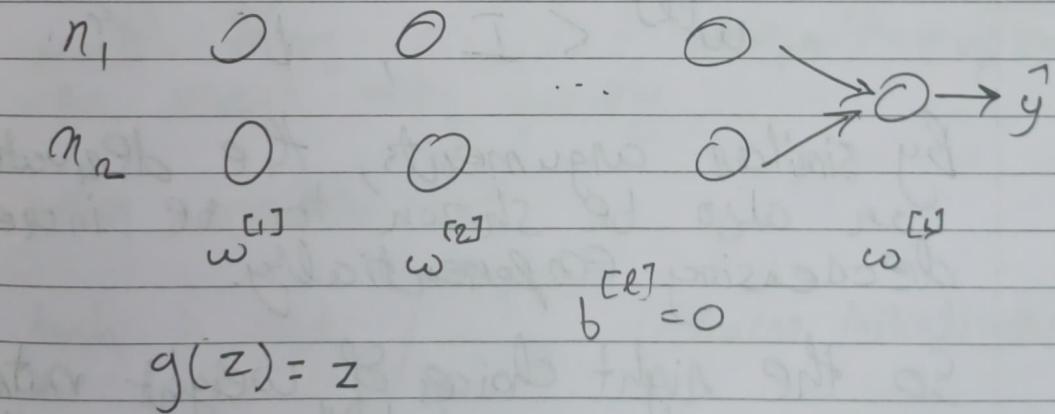
Hence normalize that.

Use the same μ and σ for the test set as well, do not calculate that separately for test set because we want the test set to go through the same transformation defined by same σ and μ calculated on training set.

Why do we Normalize ??

- Unnormalized inputs: lead to an elongated cost function and inefficient gradient descent path which makes it oscillate many times to reach the minimized loss.
- Normalized inputs: leads to a spherical cost function and efficient gradient descent path hence the algorithm quickly leads to the minimum of cost f resulting in consistent step sizes in all directions

2) Vanishing / Exploding Gradients:



$$y = w^{[1]} w^{[2]} \dots w^{[3]} w^{[4]} \underbrace{w^{[5]}}_{z^{[4]} \rightarrow a^{[1]}} a^{[2]}$$

~~Step 1: Initialize parameters~~ $a^{[3]}$ $a^{[2]}$

$$w^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \text{ where, } l = 1, 2, \dots$$

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

This is an exponential function (1.5^l), so as the NN gets deeper ($l \uparrow$), ~~to~~ the value of \hat{y} will explode.

Conversely,

$$w^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \rightarrow \hat{y} = 0.5^l \cdot x$$

In this case, the activations end up decreasing exponentially.

Note:

$$w^{[l]} > I, \uparrow \text{ exponentially}$$

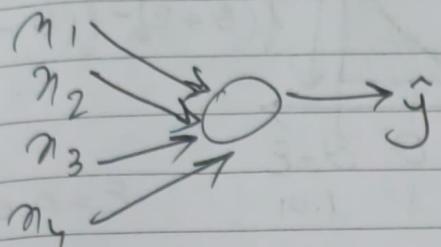
$$w^{[l]} < I, \downarrow$$

By similar arguments, the derivatives can also be shown to be increasing/decreasing exponentially.

So the right choices of weight matrices can make a lot of difference as if the gradient ~~final values~~ output gets very small then the model will take a lot of time ~~to~~ for gradient descent to learn anything as it will take small steps.

3) Weight Initialization for Deep Networks

single Neuron example



$$a = g(z)$$

$$z = w_0 n_0 + w_1 n_1 + w_2 n_2 + \dots + w_n n_n + b$$

more the # of input feature, less should be the value of w_i as they are getting summed up but if w_i is large z will become very large

$$\text{Var}(w_i) = \frac{1}{n} \quad \frac{2}{n} \text{ if activation } f \text{ is ReLU}$$

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

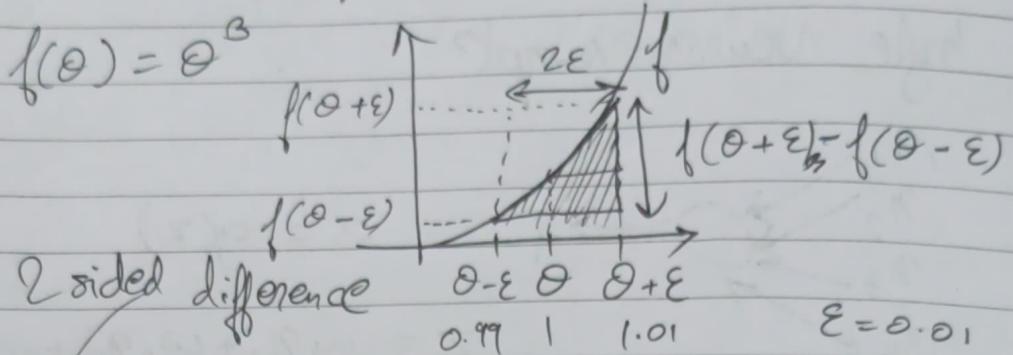
→ This will make sure that $w^{[l]}$ is not bigger than 1 and also not very small than 1 so it doesn't explode or vanish very quickly.

other variations :

$$\tanh : \text{use } \sqrt{\frac{1}{n^{[l-1]}}} : (\text{Xavier initialization})$$

$$\text{On use } \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

4) Numerical Approximation of Gradients



$$\frac{1.01^3 - (0.99)^3}{2(0.01)} = 3.001$$

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

approx error : 0.0001

5) Gradient Checking

Take $w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}$ and reshape into
in assignment: $\downarrow \downarrow \downarrow \downarrow \downarrow$ a big vector Θ .
dictionary_to_vector() \leftarrow concatenate

$$J(w^{(0)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = J(\Theta)$$

Take $d w^{(1)}, d b^{(1)}, \dots, d w^{(L)}, d b^{(L)}$ and reshape into
gradients_to_vector() $\downarrow \downarrow \downarrow \downarrow \downarrow$ a big vector $d\Theta$.

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

upto dim of

Grad check

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$$

for each i :

$$\frac{d\theta_i}{\text{approx}}^{[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \epsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

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

To check $d\theta_{\text{approx}} \approx d\theta$

$$\text{Check } \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

$$\approx 10^{-7} \rightarrow \text{correct} \\ < 10^{-3}$$

 $\|\cdot\|$: Euclidean distance $\|\cdot\|$ turns this formula into a ratiob) Implementing grad checking ~~grad~~

* Doesn't work with dropouts
 or keep-prob = 1 while working with grad check