# Ahmad Hussameldin Hamed Hassan

Shared Git-hub link: https://github.com/ahmadhassan1993/sharing-github

# Improving Deep Neural Networks: Hyperparameter tuning, Regularization and Optimization

<u>Week1 Summary:</u>

Goal is to get the suitable hybrid parameters for the specific application, by circulating this loop:
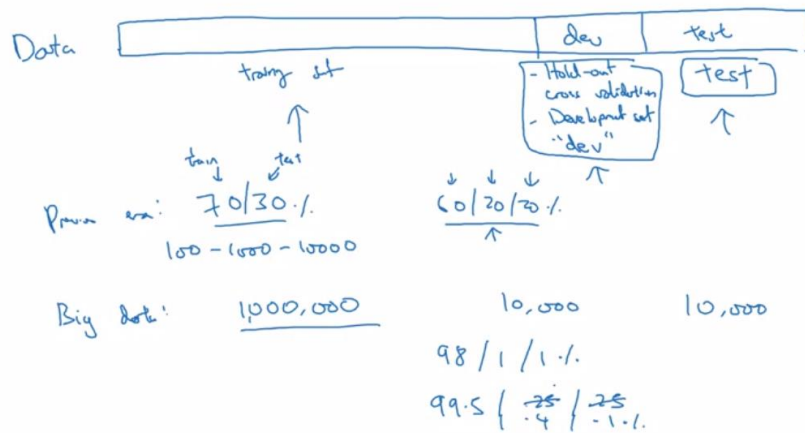


Idea: raw data from application

Code: to implement NN

Experiment: to train and test the chosen hybrid parameter. If not suitable, choose other value and do the loop again.
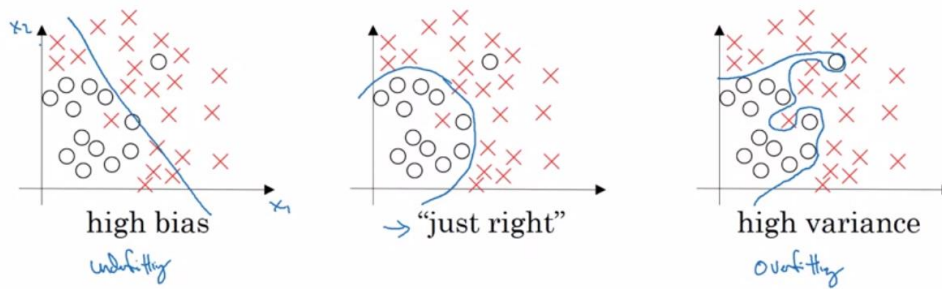
We will divide the data to three sections:

1- Training set
2- Development set
3- Test set

The second part is used to judge whether we shall change the algorithm or not (after seeing training result). This part, as well as the test part, may be large for small data set and vice versa.

Make sure to have both test or dev set and train set from same distribution, to avoid mismatch. For example, recognizing cat from website photo and personal camera photo (the resolution is different).

Not having test set is okay. We use only dev instead. Some people name dev as test, we have to be sure that it is used as dev set.

# Bias and Variance



$x_2$, $x_1$

| high bias | "just right" | high variance |
|---|---|---|
| underfitting | | overfitting |

Andrew Ng

1:49 / 8:46

---

# Bias and Variance

Cat classification

$y=1$     $y=0$



| | | | |
|---|---|---|---|
| Train set error: | 1%  15% | 15% | 0.5% |
| Dev set error: | 11%  16% | 30% | 1% |
| | high variance  high bias | high bias<br>& high varian | low bias<br>low variance |

Human $\approx 0\%$

Optimal (Bayes) error: $\approx 0\%$

Andrew Ng

## High bias and high variance



In the design:

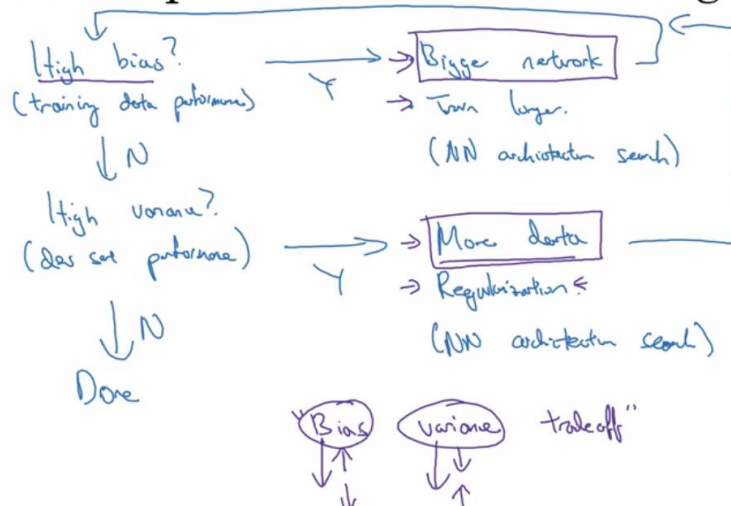If high bias=underfitting the train set, high variance=overfitting train set.

Train set error: more high, high bias

Dev set error: more high, high variance

N.B. This is based on human classifier, as the error=0

Recipe for solving both high bias and high variance:



There were a tradeoff between bias and variance, but large network and data set solves this(however with longer time).

Regularization for high variance:

## Logistic regression

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

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

lambda    lambd

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

omit

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

$L_1$ regularization    $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$    $w$ will be sparse

Andrew Ng

Lambda is a hybrid parameter modified in iterations by dev set. We use l2 regularization more than l1, as the later is sparse (have many zeros).

## Neural network

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

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

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

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

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

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

$$\frac{\partial J}{\partial w^{[l]}} = dw^{[l]}$$

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

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

Andrew Ng

We use Forbenius norm regularization in NN instead of L2 in logistic regression. This will lead to Weight Decay formula in W after back propagation.
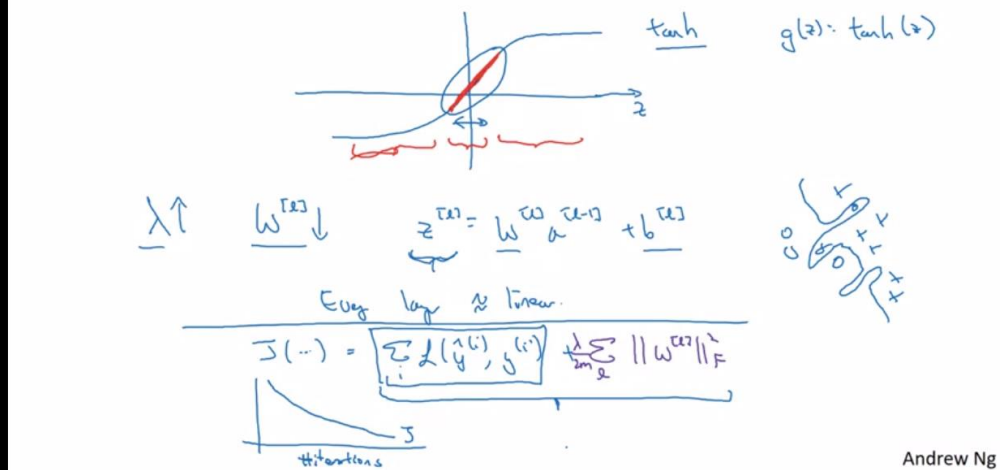
How does Regularization reduces overfitting (high variance):

How does regularization prevent overfitting?
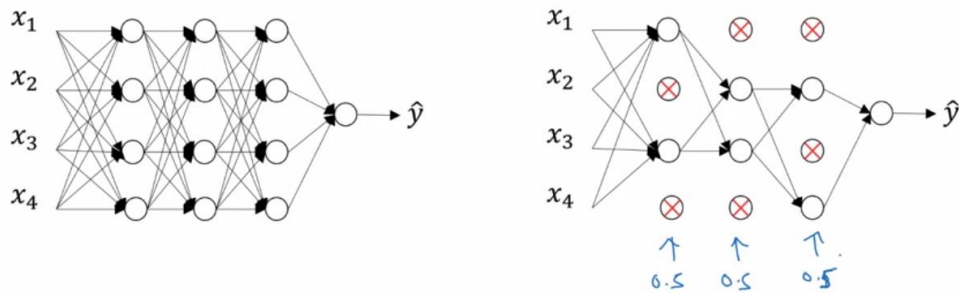
Higher lambda, lower W (which is not practical). So, lower Z, which tend to make activation functions linear. This is not practical for Deep NN at all.



How does regularization prevent overfitting?

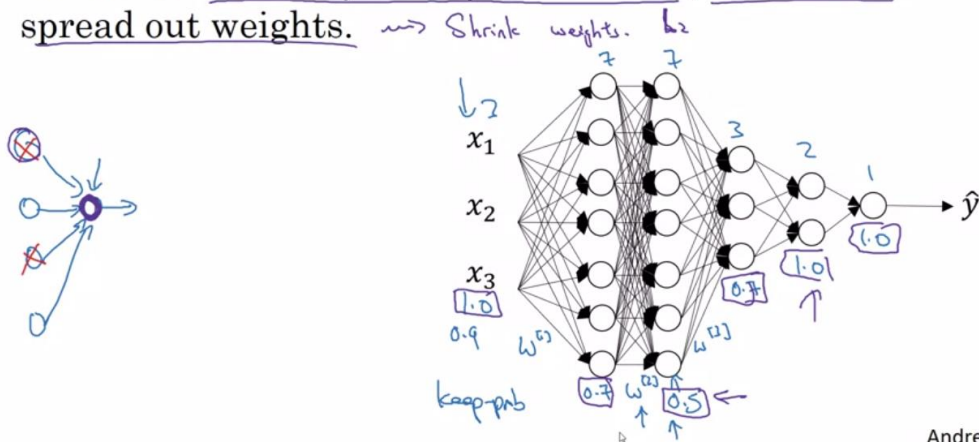Drop-out regularization:

For forward propagation, Inverted Drop-out:



Famous type is Inverted Reg. We will use also d3 for back propagation.

Drop-out is not implemented in test time, as this will add noise.

Why does drop-out work?

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

The drop-out is highly used in larger dimensions of hidden layers, and not used for both inputs and output layers (for importance of each element existence).

Downside: cost function is hard to be defined.

Two methods are now used to reduce both cost function and overfitting:



1- Early stopping: this is an orothogonalization method, where we use gradient descent to reduce cost function and regularization to reduce overfitting.
2- Using l2 regularization with large iterations for lambda.

Normalizing inputs will make the cost function diverge faster:

# Normalizing training sets

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



Subtract mean:

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

$$X := x - \mu$$

Normalize variance

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

$$X /= \sigma^2$$

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

2:15 / 5:30    Andrew Ng

# Why normalize inputs?

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

$w_1, x_1 : 1 \cdots 1000 \leftarrow$
$w_2, x_2 : 0 \cdots 1 \leftarrow$
$-1 \cdots 1$

Unnormalized:

Normalized:



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

4:57 / 5:30    Andrew Ng

Gradients can increase or decrease exponentially:

Andrew Ng

To solve this, we initialize the weights such that it will not be more (greater or less) than 1. This depends on the type of the activation function:



To check the gradient computation, we have to do the numerical method:

## Checking your derivative computation

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

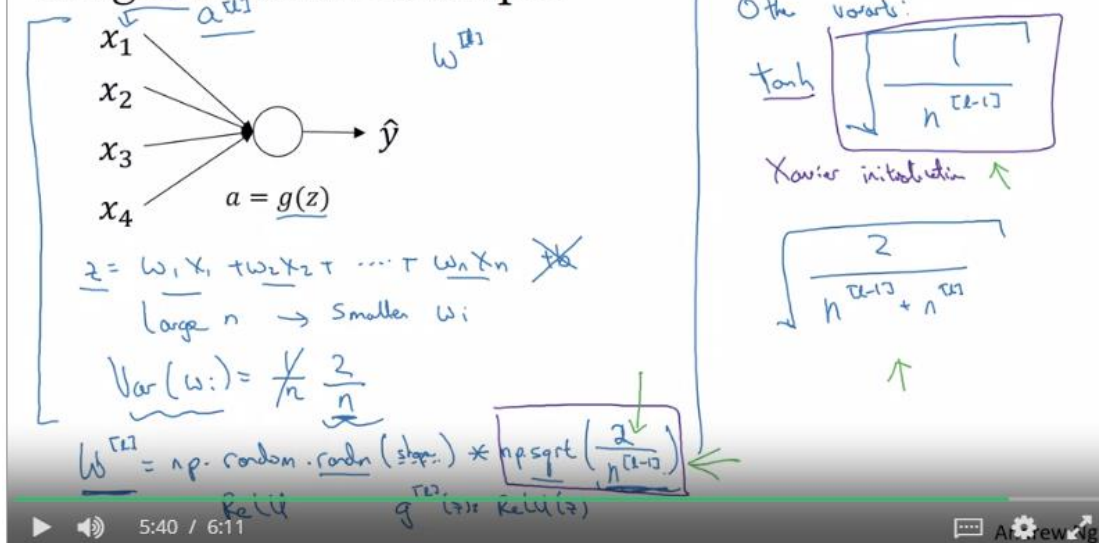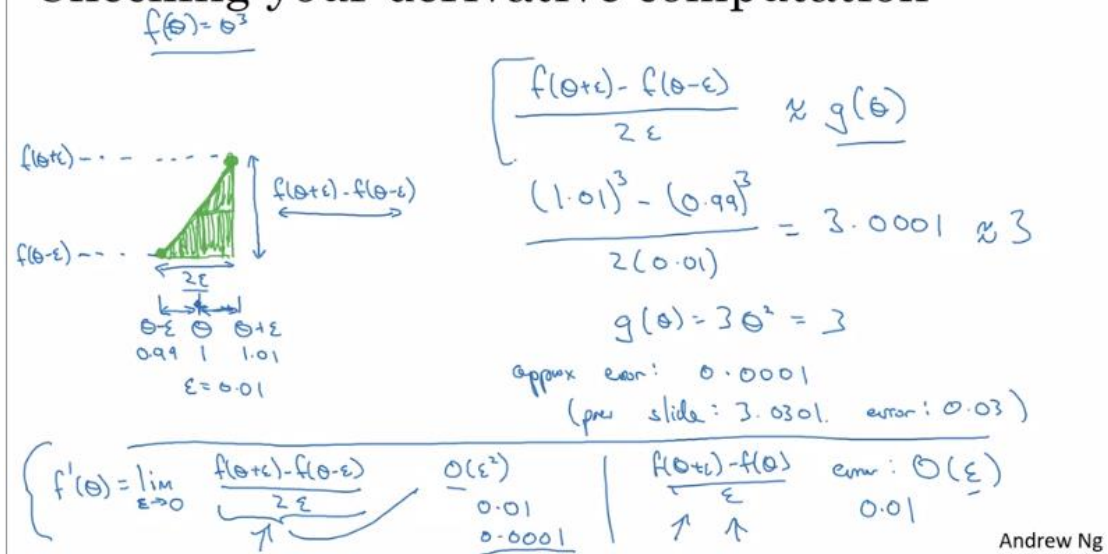$$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \approx g(\theta)$$

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

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

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

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

0.01
0.0001

0.01

$f(\theta+\varepsilon)$
$f(\theta-\varepsilon)$
$f(\theta+\varepsilon) - f(\theta-\varepsilon)$
$2\varepsilon$
$\theta-\varepsilon \quad \theta \quad \theta+\varepsilon$
$0.99 \quad 1 \quad 1.01$
$\varepsilon = 0.01$

First, we make a large variable Θ:

## Gradient check for a neural network

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

concatenate

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

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

concatenate

Is $d\theta$ the gradient of $J$

1:45 / 6:34

Then do the check procedure as follows:

Gradient checking (Grad check)

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

for each $i$:

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

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

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

$\varepsilon = 10^{-7}$

$\approx \boxed{10^{-7} \quad - \text{great!}} \leftarrow$

$10^{-5}$

$\to 10^{-3} \quad - \text{worry.} \leftarrow$

6:16 / 6:34                                                                 Andrew Ng

If the check equation equals Ɛ or very near, so it is ok. Otherwise, try to debug the Ө(i) variables.

Important notes on applying gradient descent:



Gradient checking implementation notes

- Don't use in training – only to debug   $d\theta_{approx}[i] \longleftrightarrow \dfrac{d\theta[i]}{}$

- If algorithm fails grad check, look at components to try to identify bug.
  $db_k^{[l]} \qquad dw_k^{[l]}$

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

  $d\theta = \text{gradt of } J \text{ w.r.t. } \theta$

- Doesn't work with dropout.   $J$   keep-prob $= 1.0$

- Run at random initialization; perhaps again after some training.
  $w, b \approx 0$

Andrew Ng

N.B. for drop-out regularization: we first run the gradient check without drop-out, then apply the drop-out.