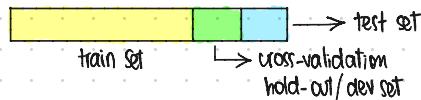


# HYPERPARAMETER TUNING, REGULARIZATION AND OPTIMIZATION

COURSE II

## PRACTICAL ASPECTS OF DEEP LEARNING

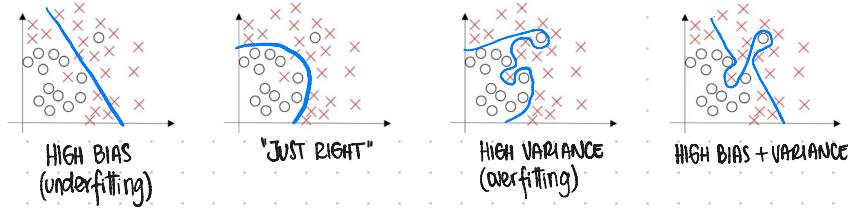
- #1. train / developer / test sets:



Before: 70/30% - 60/20/20%

Big data era: 98/1/1% - 99.5/0.4/0.1%

- #2. bias vs. variance: is there a trade-off?



train set error: 1%

15%

15%

0.5%

dev. set error: 11%

16%

30%

1%

HIGH VAR.

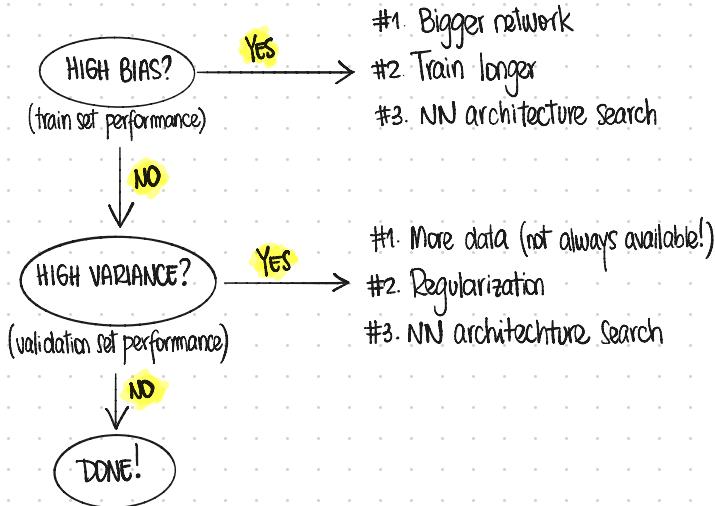
HIGH BIAS

HIGH BIAS  
+ HIGH VAR.

LOW BIAS  
+ LOW VAR.

WORST NIGHTMARE

## BASIC RECIPE FOR MACHINE LEARNING



## REGULARIZATION

- \* prevent/reduce variance (overfitting)

- \* regularization in logistic regression:

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

REGULARIZATION PARAMETER

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

L<sub>2</sub> REGULARIZATION

(most used method)

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

NOT NECESSARY  
( $b$  is just one parameter  
over a large # of parameters)

$$\min J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_1$$

L<sub>1</sub> REGULARIZATION

$$\|w\|_1 = \sum_{j=1}^n |w_j|$$

L<sub>1</sub> REGULARIZATION

w will be sparse (with lots of 0)  
not used very often

\* regularization in neural networks:

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

FROBENIUS NORM  
( $L_2$  for matrices)

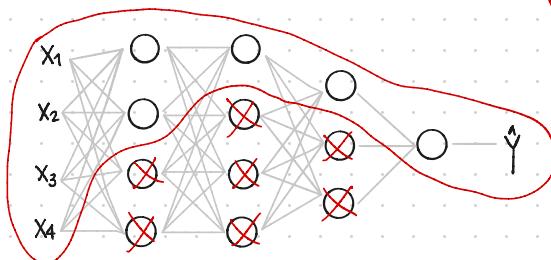
$$\|w^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l+1]}} (w_{ij}^{[l]})^2$$

$w(n^{[l]} \times n^{[l+1]})$

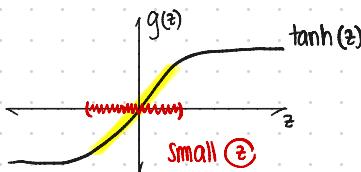
WEIGHT DECAY

\* why does regularization prevent overfitting?

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



WHEN WE CHOOSE A LARGE  $\lambda$ , WE MAKE  $w^{[l]} \approx 0$ . THIS REDUCES THE IMPACT OF THE HIDDEN UNITS IN THE NETWORK, "MAKING IT SMALLER". SINCE THE NETWORK IS NOW SMALLER, IT WILL BE ABLE TO GENERALIZE BETTER.

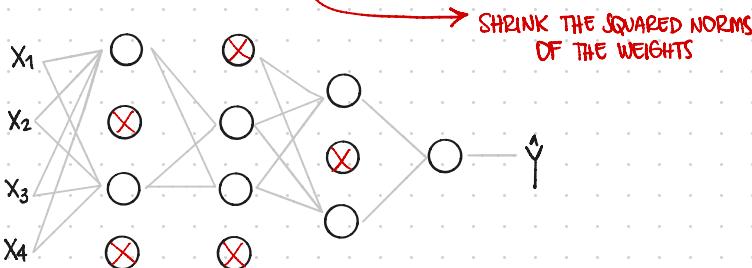


\* when  $\lambda$  is large,  $w^{[l]}$  will be small.

Because  $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$  will be small,  $g(z)$  will become linear. Thus, every layer will be roughly linear as in linear regression. Even a very deep network will have these "linear" characteristics and will do a much better job preventing overfitting.

## DROPOUT REGULARIZATION

- \* it consists of cancelling some hidden units and make the network "smaller"
- \* it works, because given that the units can't rely on any one feature (it may get knocked out), it has to spread out the weights.



- \* downside  $\rightarrow$  the cost function  $J(w, b)$  is no longer well-defined

## OTHER REGULARIZATION METHODS

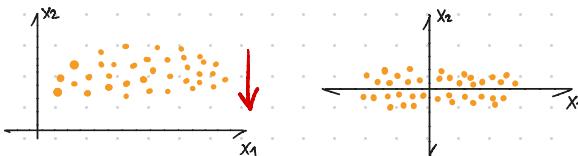
- #1. data augmentation
- #2. early stopping  $\rightarrow$  Stop training your network when dev set error  $>$  train set error. By stopping half-way, you end up with a mid-size  $\|w\|_F^2$ .
  - downside: orthogonalization. By stopping the training you stop the optimization and also cannot tackle overfitting separately.

## NORMALIZING INPUTS

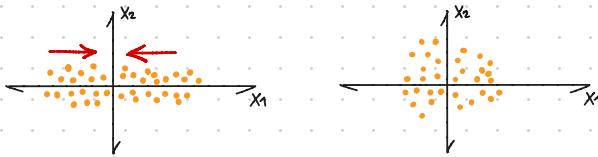
- \* speeds up your training

- \* steps:

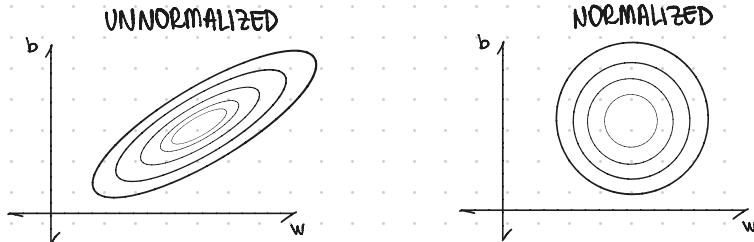
(#1). Subtract the mean  $\rightarrow \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \Rightarrow x = x - \mu$



#2. normalize variance  $\rightarrow \sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)2} \Rightarrow x / \sigma^2$



\* you need to transform both the train/test set.

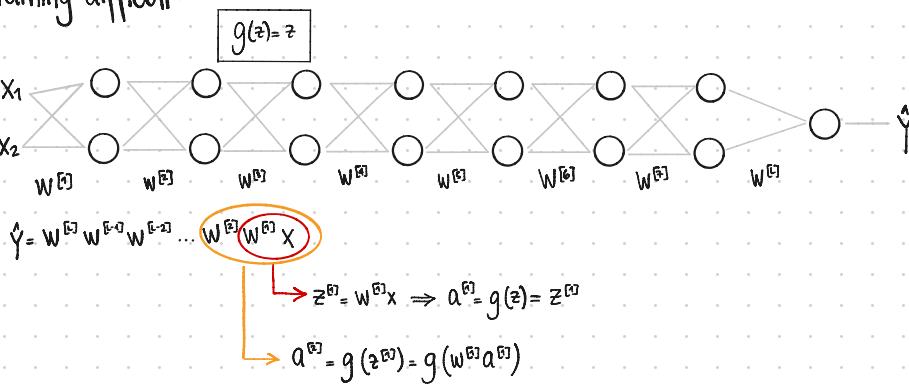


\* the gradient descent needs to take more steps to the min.

\* the gradient descent takes less steps towards the minimum.

## VANISHING/EXPLODING GRADIENTS

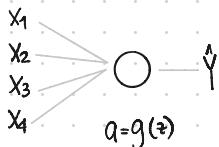
\* when you train very deep networks, your slope can get either very big or small, making training difficult.



If  $W[1] = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ ,  $\hat{y} = W[L] \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \approx 1.5^{L-1} \hat{y}$  WILL GROW EXPONENTIALLY AND WILL EXPLODE!!

If  $W^{[1]} = \begin{bmatrix} 0.5 & & \\ 1.5 & 0 & \\ 0 & 1.5 & \\ & 0.5 & \end{bmatrix}$ ,  $\hat{y} = W^{[1]} \begin{bmatrix} 0.5 & & \\ 1.5 & 0 & \\ 0 & 1.5 & \\ & 0.5 & \end{bmatrix}^{L-1} X \approx 1.5^{L-1} \rightarrow \hat{y}$  WILL DECREASE EXPONENTIALLY AND WILL VANISH!!

## WEIGHT INITIALIZATION FOR DEEP NETWORKS



$$z = W_1 x_1 + W_2 x_2 + \dots + W_n x_n + b$$

LARGER ( $n$ ), SMALLER ( $W_i$ )

$$\text{Var}(W_i) = \frac{1}{n} \Rightarrow W^{[1]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{(L-1)}}\right)$$

HYPERPARAMETER

$$\text{Var}(W_i) = \frac{2}{n} \Rightarrow \text{ReLU activation}$$

$$\text{Var}(W_i) = \sqrt{\frac{1}{n^{(L-1)}}} \Rightarrow \text{tanh activation (Xavier activation)}$$

## GRADIENT CHECKING

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

CONCATENATE so...

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

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

\* for each  $i$ :

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

\* check:  $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta\|_2} \approx 10^{-7}$  ✓ GOOD!

$\frac{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}{\|d\theta\|_2} \approx 10^{-3}$  ✗ WORRY!

\* implementation notes...

- #1. don't use in training - only to debug
- #2. if algorithm fails the grad check, look at the components to try identify a bug.
- #3. remember regularization!
- #4. doesn't work with dropout (implement grad check without dropout and then do dropout).
- #5. run at random initialization

## MINI-BATCH GRADIENT DESCENT

\* split the training set into mini-batches

$$X = \left[ \underbrace{x^{(1)} \ x^{(2)} \ x^{(3)} \ x^{(4)} \dots \ x^{(1000)}}_{X^{[t]} \ (n, 1000)} \mid \underbrace{x^{(1001)} \dots \ x^{(2000)}}_{X^{[t]} \ (n, 1000)} \mid \dots \ x^{(n)} \right] \quad \left. \right\} \text{mini-batch } t: X^{[t]}, y^{[t]}$$
$$Y = \left[ \underbrace{y^{(1)} \ y^{(2)} \ y^{(3)} \ y^{(4)} \dots \ y^{(1000)}}_{Y^{[t]} \ (1, 1000)} \mid \underbrace{y^{(1001)} \dots \ y^{(2000)}}_{Y^{[t]} \ (1, 1000)} \mid \dots \ y^{(n)} \right]$$

for  $t = 1, \dots, 5000$ :  $\rightarrow$  # MINI-BATCHES

forward propagation on  $X^{[t]}$ :

$$Z^{[t]} = W^{[t]} X^{[t]} + b^{[t]}$$

$$A^{[t]} = g^{[t]}(Z^{[t]})$$

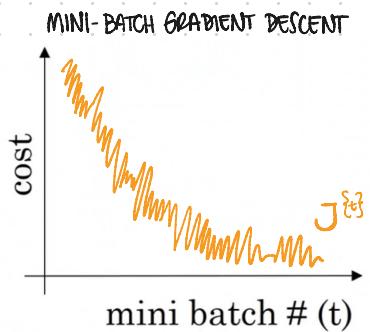
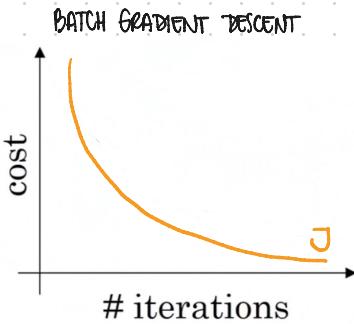
$$\vdots$$
  
$$A^{[t]} = g^{[t]}(Z^{[t]})$$

Compute cost  $J^{[t]} = \frac{1}{1000} \sum_{i=1}^l \mathcal{L}(y^{(i)}, A^{[t]})$

backprop.  $\rightarrow$  # EXAMPLES IN THE MINI BATCH

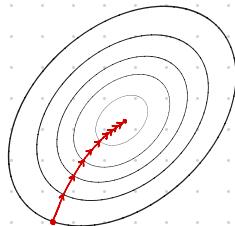
EXAMPLES FROM  $X^{[t]}, Y^{[t]}$

\* it runs much faster!



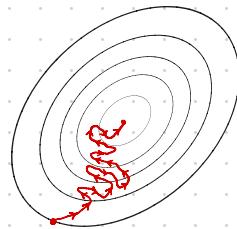
if mini-batch size =  $m$  # TRAINING SAMPLES

**BATCH GRADIENT DESCENT**

$$(X^{[1]}, Y^{[1]}) = (X, Y)$$


\* too long per epoch

if mini-batch size = 1  
**STOCHASTIC GRADIENT DESCENT**  
every example is its own mini-batch



\* lose all the speedup from vectorization

WHAT WORKS BEST?



SOMETHING IN BETWEEN!

mini-batch size not  
too big / too small  
↓  
fastest learning

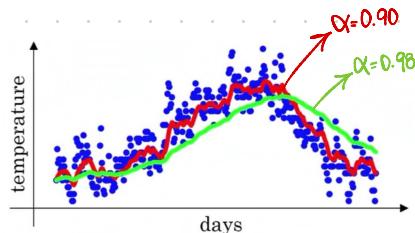
typical mini-batch size:

64, 128, 256, 512  
 $\downarrow$   
 $2^6$     $2^7$     $2^8$     $2^9$

## EXPONENTIALLY WEIGHTED AVERAGES

$$V_t = \alpha V_{t-1} + (1-\alpha) v_t$$

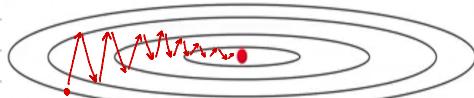
CONSTANT  
MOVING AVERAGE



- \* if  $\alpha$  is close to 1, it gives more importance to  $V_{t-1}$
- \* a smaller  $\alpha$  is more susceptible to outliers.

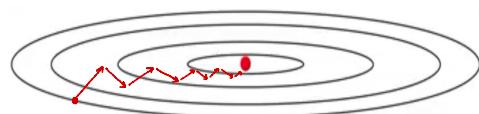
## GRADIENT DESCENT WITH MOMENTUM

- \* compute exponentially weighted averages, and use these averages to compute gradient descent



PLAIN GRADIENT DESCENT

\* here, you want a slower learning along the vertical axis to prevent overshooting. You also want a faster learning across the horizontal axis. So...



GRADIENT DESCENT + MOMENTUM

compute  $dW, db$  on current mini-batch

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

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

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

MOMENTUM TERMS

## ROOT-MEAN-SQUARED PROP

- \* compute  $dW, db$  on current mini-batch.

SMALL  $S_{dw} = \beta S_{dw} + (1-\beta) dw^2$

LARGE  $S_{db} = \beta S_{db} + (1-\beta) db^2$

$$W = W - \alpha \frac{dw}{\sqrt{S_{dw}}}, \quad b = b - \alpha \frac{db}{\sqrt{S_{db}}}$$

- \* these derivatives are much larger in the vertical direction than in the horizontal direction

