

## Training a neural network

- 1) Randomly initialize weights ( $\approx 0$ )
- 2) Implement forward propagation to get  $h_\theta(x^{(i)})$  for any  $x^{(i)}$
- 3) Implement code to compute cost function  $J(\theta)$
- 4) Implement backprop to compute partial derivatives  $\frac{\partial J(\theta)}{\partial \theta_{jk}^{(l)}}$

for  $i = 1:m$

perform forward prop & back prop using  $(x^{(i)}, y^{(i)})$

(get activations  $a^{(l)}$  &  $z^{(l)}$  for  $l = 2, \dots, L$ )

5) Use gradient checking to compute  $\frac{\partial}{\partial \theta_{j,n}} J(\theta)$  computed using

back propagation vs using numerical estimates of gradient of  $J(\theta)$   
Then disable gradient checking code.

6) Use 'gradient descent or advanced optimization method with back propagation to try to minimize  $J(\theta)$  as a function of parameters  $\theta$ .

## How to debug a learning algorithm

- Get more training example - fixes high variance.
- Try smaller set of features - fixes high variance.
- Try getting additional features - fixes high bias.
- Try adding polynomial features - fixes high bias.  
(fixes high bias) ( $x_1^2, x_2^2, x_1 x_2$ , etc)
- Try decreasing  $\lambda$  (Regularized expression)
- Try increasing  $\lambda$ . (fixes high variance)

## Machine learning diagnostic

: A tests you can run to

gain insight about the learning algorithms

→ 70% train

→ 30% test

To do,

→ learn parameter  $\theta$  from training data  
(minimizing training error)  $J(\theta)$ .

→ compute test set error. (linear regression)

$$J_{\text{test}}(\theta) = \frac{1}{2m} \sum_{i=1}^{m_{\text{test}}} [h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)}]^2$$

Model selection

60% train set

20% cross validation set

20% test set.

\* when we have a lot of features, we use regularization ~~for overfit~~

high bias  $\rightarrow \lambda$  is large  $\therefore$  all other features of hypothesis are  $\approx 0$   
 $\therefore h_0(x) \approx 0$

high variance  $\rightarrow \lambda \approx 0$ .

Learning curves

high bias won't help if you increase no. of examples

high variance, get more data

Building a spam classifier,

$x$  = features of email.

$y$  = spam(1) or not spam(0).

$x$  = look for words that are used in spam i.e., deal, buy, now, discounts. (or Andrew if we use it it's not a spam).

How to go about it

- Collect data
- Develop algorithm based on the email routing (contained in the header)
- Develop sophisticated algo - i.e., 'discounts' & 'discounted' are same or not, 'deal' & 'deals' are same or not.
- Develop sophisticated algorithm to detect misspellings (watches)

### Approach

- Start simple
- plot learning curve & decide whether more data or more features.
- Error analysis - Manually examine the examples

Numerical evaluation - Single number for checking the error.

Skewed class :- If you have only examples which point to one class i.e., patient doesn't have cancer.

		Actual class	
		1	0
Predicted class	1	True positive	False positive
	0	False negative	True Negative

$$\frac{85}{85+890}$$

Precision - of all patients predicted, who actually have cancer

$$\text{precision} = \frac{\text{True pos.}}{\text{True pos.} + \text{False pos.}}$$

Recall - of all patients that actually have cancer, how much did we predict

$$\text{recall} = \frac{\text{True pos.}}{\text{True pos.} + \text{False neg.}}$$

for precision recall ~~weights~~ trade off

$$F_1 \text{ score} = \frac{2PR}{P+R}$$