

## Week 6

If when you test your hypothesis you get a high accuracy, but it fails on new data, try the following:

- (1) \*
- Get more training ~~ex~~ examples
  - Try smaller sets of features
  - Try adding polynomial features
  - Try adding additional features
  - Try decreasing  $\lambda$
  - Try increasing  $\lambda$
- This is used later in the notes

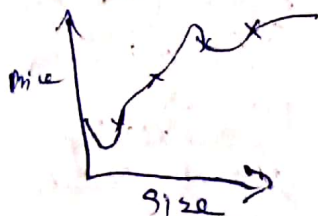
How do you decide which to do?

Machine Learning diagnostic:

Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

## Evaluating a hypothesis

If overfitting occurs, then it may show high accuracy but fail on new data



On small no. of features we can plot and see, but what if we have many features?

We split the data: notation

Training set - 70%  $m$

Test set - 30%  $m_{\text{test}}$   $(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$

Distribute them randomly (first randomly sort the data, then distribute)

Training / Testing Procedure (for linear regression)

- learn parameter  $\theta$  from training data  
(Minimising training error  $J(\theta)$ )

- compute test set error

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

For logistic regression,

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$$

- Misclassification error

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5, y=0 \\ & h_{\theta}(x) \leq 0.5, y=1 \\ 0 & \text{otherwise} \end{cases} \quad \text{error}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

## Model selection

How to choose polynomial features,  
d - degree of polynomial

d=1 1.  $h_0(x) = \theta_0 + \theta_1 x \xrightarrow{\text{min } J(\theta)} \theta^{(1)}$  - This is the minimised cost function's  $\theta$

d=2 2.  $h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)}$

d=3 3.  $h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)}$

d=10 10.  $h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)}$

Find  $\theta$  for the different models

Then  $J_{\text{train}}(\theta^{(1)})$ ,  $J_{\text{train}}(\theta^{(2)})$ ...  $J_{\text{train}}(\theta^{(3)})$

and we choose which has the least cost

But note that it fit well to this training test, and may not fit well for new data

$\therefore$  We split the data as follows

- Training set - 60% -  $m$
- Cross Validation Set - 20% -  $m_{cv}$   $(x_{cv}^{(i)}, y_{cv}^{(i)})$
- Test set - 20% -  $m_{\text{test}}$   $(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$

Train/validation/test error

Training error:

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

Cross validation error:

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_0(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_0(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$



Now we choose model as follows:

$$1. h_0(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$$

$$2. h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$$

$$3. h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$$

$$10. h_0(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})$$

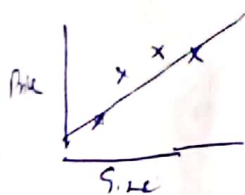
So we fit the  $\theta$  on the cross validation test and pick the  $\theta$  best model.

Let's say 4 was the best,

Now we estimate generalization error for test set  $J_{test}(\theta^{(4)})$

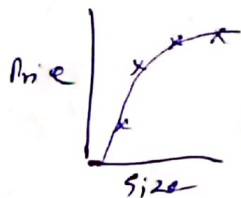
Diagnosing bias vs variance:

Bias / Variance



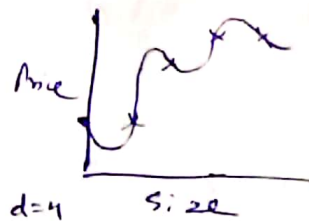
$$d=1 \quad \theta_0 + \theta_1 x$$

High Bias  
Results in underfit



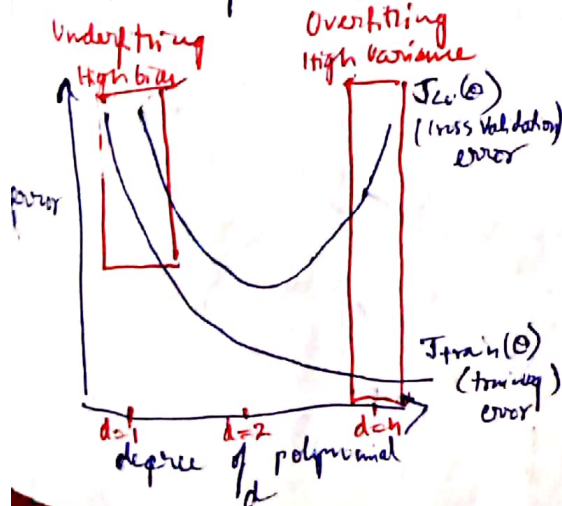
$$d=2 \quad \theta_0 + \theta_1 x + \theta_2 x^2$$

Just right



$$d=4 \quad \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance  
Results in overfit



Bias (underfit):

$\rightarrow J_{train}(\theta)$  will be high

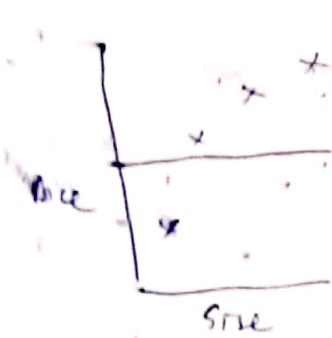
$\rightarrow J_{cv}(\theta)$  will be high ( $J_{cv}(\theta) \approx J_{train}(\theta)$ )

Variance (overfit):

$\rightarrow J_{train}$  will be low

$\rightarrow J_{cv}(\theta) \gg J_{train}(\theta)$

# Regularisation and bias/variance

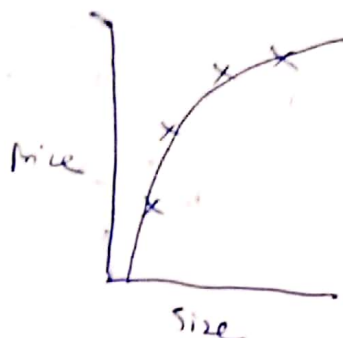


Large  $\lambda$

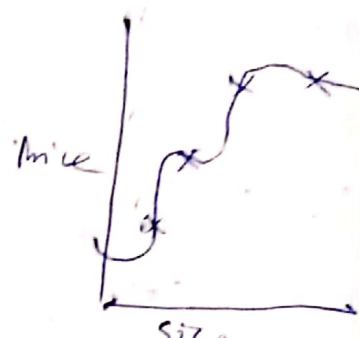
High bias (underfit)

$\lambda = 10000, \theta_1 \approx 0, \theta_2 \approx 0 \dots$

$h_0(x) \approx \theta_0$



Intermediate  $\lambda$



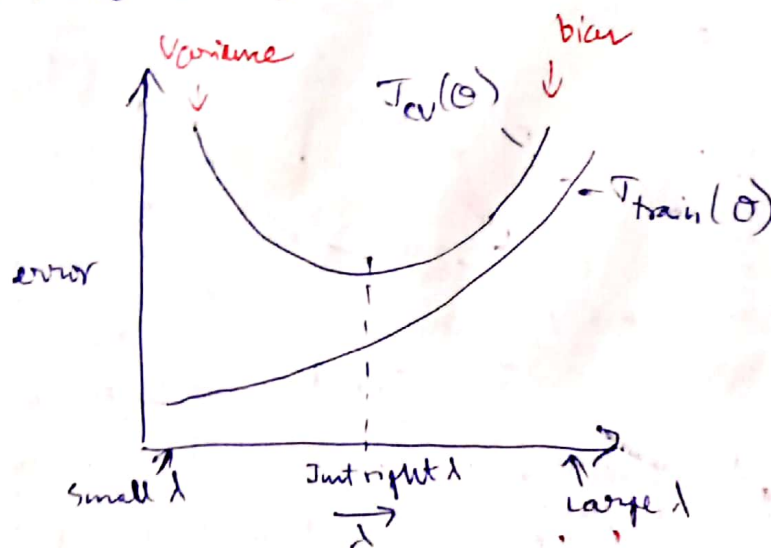
Small  $\lambda$

High variance (overfit)

$\lambda = 0$

Crossing  $\lambda$ :

1. Try  $\lambda = 0 \xrightarrow{\min J(\theta)} \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
  2. Try  $\lambda = 0.01 \xrightarrow{\min J(\theta)} \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
  3. Try  $\lambda = 0.02 \xrightarrow{\min J(\theta)} \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
  4. Try  $\lambda = 0.04 \xrightarrow{\min J(\theta)} \theta^{(4)} \rightarrow J_{cv}(\theta^{(4)})$
  - $\vdots$
  12. try  $\lambda = 10 \xrightarrow{\min J(\theta)} \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
- Pick the best  $\theta$  (lets say  $\theta^{(5)}$ )  
Then  $J_{test}(\theta^{(5)})$



$$h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{i=1}^m \theta_i^2$$

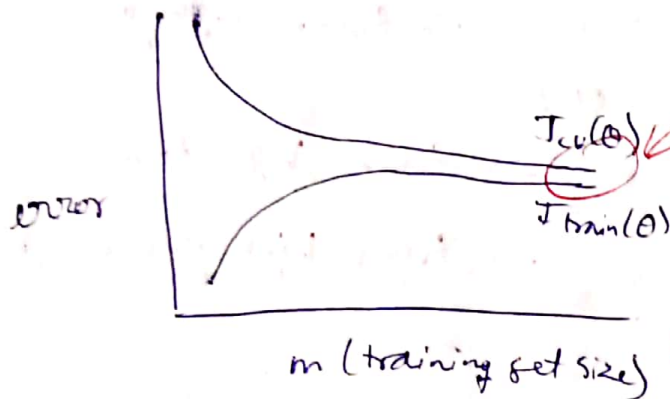
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

$$J_{cv}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

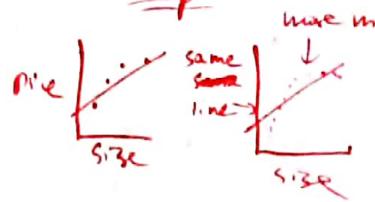
## Learning curve

its a tool to find if the model is suffering from underfitting or overfitting

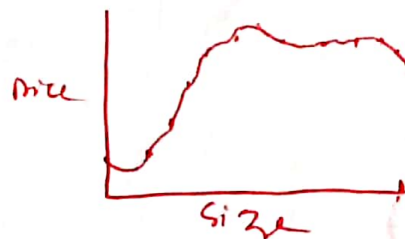
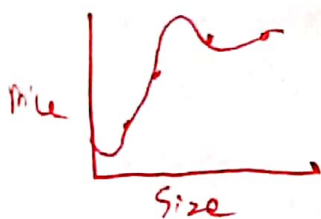
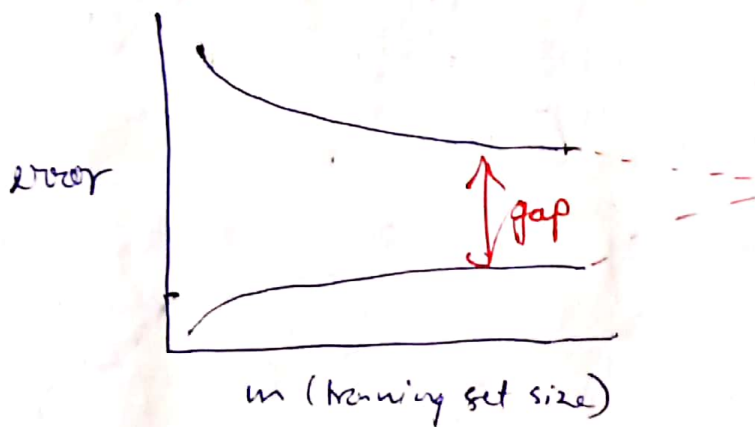
Experiencing high bias:



After certain  $m$ , the graph flattens out, so getting more data won't help



Experiencing high variance



$\therefore$  GETTING MORE TRAINING  
EXAMPLES (INCREASING  $m$ )  
HELPS ONLY FOR HIGH VARIANCE



## Looking back at ①\*

- Get More training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features → fixes high bias  
( $x_1^2, x_2^2, x_1 x_2$ , etc)
- try decreasing  $\lambda$  → fixes high bias
- Try increasing  $\lambda$  → fixes high variance

## Neural Networks:

### Small NN



- Few parameters
- Computationally cheaper
- prone to underfitting

~~No~~

### Large NN



→ More units



→ More layers

- Computationally more expensive
- Prone to overfitting
- Use regularization

To decide no. of layers:

- Plot  $J_{cv}(l)$  with different no. of layers and see which fits best

# Building a Spam classifier

Supervised learning.  $x$  = features of email  
 $y$  = spam (1) or not spam (0)

Features  $x$ : Choose 100 words indicative of spam/not spam (eg. deal, buy, discount, andrew, now, ...)

$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ \vdots \end{bmatrix}$  andrew  
buy  
deal  
discount  
now  
:

From: [deepsales@buystufffromme.com](mailto:deepsales@buystufffromme.com)  
To: [cong@cs.stanford.edu](mailto:cong@cs.stanford.edu)  
Subject: Buy now!  
Deal of the week! Buy now!

Indicates if the word appeared in the email

Not spam example

Note: In practice, take most frequently occurring  $n$  words (10,000 to 50,000) in training set, rather than manually pick 100 words.

How to spend how time to make it have low error?

- Collect lots of data
  - eg. Honey pot project
- Develop sophisticated features based on email routing information (from email header)
- Develop sophisticated features for message body, eg. should "discount" and "discounts" be treated as the same word? How about "deal" and "dealer"?
- Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (eg. mortgage, medicine, withder).



## Error Analysis

Recommended approach:

- Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross-validation data.
- Plot learning curves to decide if more data, more features, etc are likely to help.
- Error analysis: Manually examine the errors on examples in the cross validation set and try to spot a trend where most of the errors were made.

The importance of numerical evaluation

It is important to get error results as a single, numerical value, in order to assess the algorithm's performance.

For example, in the spam classifier, if we use stemming (treating similar words like discount/discounts/discounted as same word) we get a error rate of 3% instead of 5%. Without it, then we can decide to add it to the model. So using numerical evaluation we can determine whether to keep a feature or not.

Error metrics for skewed classes:

Consider a cancer classification example. And let's say we get an error of 1%. These seem great since accuracy is 99%. However let's say only 0.50%  $\left(\frac{0.5}{100}\right)$  patients have cancer.

Then a non-learning algorithm that always returns  $y=0$  (person has no cancer) will have an error of ~~0.5%~~ 0.5%.

So a skewed class is when one class is substantially lower in number than the other class such that an algorithm that always returns the majority class will have less error.

$\therefore$  We need another evaluation metric

### Precision/Recall

$y=1$  is presence of rare class that we want to detect. Actual class

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

Precision (what fraction actually have cancer?  
of all the ppl we predicted  $y=1$ )

$$\frac{\text{True Positive}}{\text{\# predicted positives}} = \frac{\text{True +ve}}{\text{True +ve} + \text{False +ve}}$$

$\therefore$  A model that always predicts  $y=0$ , will have recall = 0

Recall (what fraction did we correctly detect as having cancer of all the patients that actually have cancer)

$$\frac{\text{True +ve}}{\text{\# actual true}} = \frac{\text{True +ve}}{\text{True +ve} + \text{False -ve}}$$

## Trading off precision and recall

Logistic regression:  $0 \leq h_0(x) \leq 1$

Predict 1 if  $h_0(x) \geq \text{threshold}$

Predict 0 if  $h_0(x) < \text{threshold}$

- usually we used to set threshold as 0.5

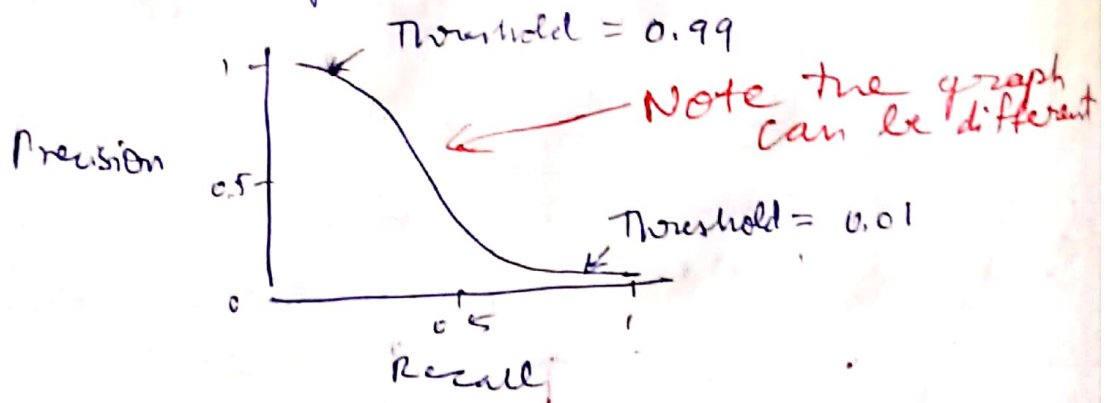
Consider two scenarios:

1) Suppose we want to predict  $y=1$  (cancer) only if we are very confident. Then we set threshold = 0.7 or 0.9 or some high number.

→ higher precision, lower recall

2) Suppose we want to avoid missing too many cases of cancer (avoid false negatives), then we use a low threshold like 0.3.

→ higher recall, lower precision



Can we choose the threshold automatically?

Before we had 1 numerical evaluation but now we have 2, which makes choosing an algorithm from another hard.



Using average isn't a good idea.

F<sub>1</sub> score:  $2 \frac{P R}{P + R}$  ← F score method

	Precision (P)	Recall (R)	F <sub>1</sub> Score
Algo 1	0.5	0.4	0.444
Algo 2	0.7	0.1	0.175
Algo 3	0.02	1.0	0.0392

Data for ML

- learning
- 1) We use a low-bias algorithm (logistic/linear regression with many <sup>units</sup> features, NN with many hidden ~~layers~~). Due to the large number of parameters it'll have low bias.  $J_{\text{train}}(\theta)$  will be small.
  - 2) Use a large training set (~~undoubtedly~~ ~~to overfit~~) so overfitting is unlikely to happen.

$$J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$$

∴  $J_{\text{test}}(\theta)$  will be small.

We need to also have enough features such that a human expert can predict y.  
Exa: Just given a house size, a human expert can't predict the price (since locality, furnished or not, etc matters). So we need to add all these features.