

# Evaluating a learning algorithm

What happens when our learning algorithm does not predict well?  
What can we do?

- You can try
- Get more training data
  - try smaller sets of features
  - Try getting additional features
  - Try with polynomial features
  - Try increasing/decreasing  $\lambda$

But which one?  
Intuition?

Better run

diagnostic tests  
and get insights  
about Learning Alg.  
performance so we  
know how to improve  
it.

## Evaluate a hypothesis

- Split the dataset into

training set  $(\approx 70\%)$

test set  $(\approx 30\%)$

random selection  
is better

- Training/testing for linear regression

- Learn  $\Theta$  from training data minimizing  $J(\Theta)$  (from 70% training data)

$$\text{Compute test set error} \rightarrow J_{\text{test}}(\Theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\Theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

- Training/testing for logistic regression

- Learn  $\Theta$  from training data minimizing  $J(\Theta)$  (from 70% training data)

$$\text{Compute test set error} \rightarrow J_{\text{test}}(\Theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\Theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log (1 - h_{\Theta}(x_{\text{test}}^{(i)}))$$

or

Compute missclassification error

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

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

average of number  
of misclassifications

## Model selection with train/cross-validation/test sets

We might want to use different models (e.g. different polynomial degrees). Each model gives a set of  $\Theta$ , which one do we choose?

↓  
We can select the model that has a lower error in the test set

60% train ← train and get  $\Theta$  for all models (fit  $\Theta$ )  
20% cross-validation ← choose the best model here! (fit  $\Theta$ )  
20% test ← then finally test the chosen model

A

but we are using the test set  
on our learning process!

actually to fit the parameter  $\lambda$  (poly. degree)

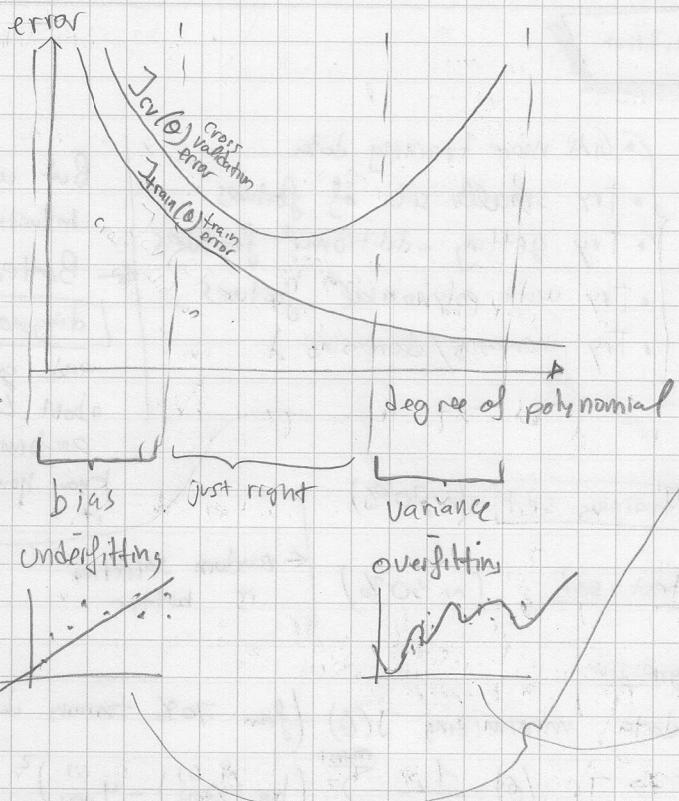


## Diagnose Bias (underfit) vs. Variance (overfit)

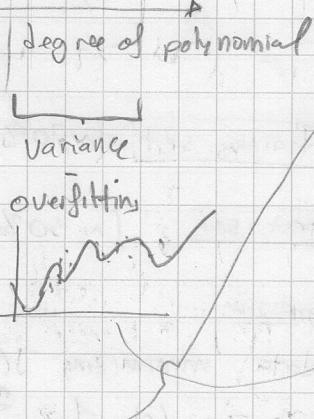
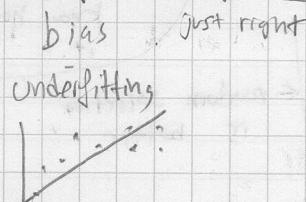
what does this mean?

bias  $\leftrightarrow$  underfit ?

Variance  $\leftrightarrow$  overfit



would be nice to see how  
under/over-fitting a model is, but if  
we have more than 1-2 features...  
we cannot plot it!



If our learning algorithm is not working well ( $J_{\text{train}}$  or  $J_{\text{cv}}$  are high), are we suffering from bias (underfit) or variance (overfit) problem?

$J_{\text{train}} \& J_{\text{cv}}$  high  $\Rightarrow$  high bias (underfit)

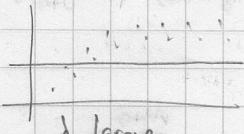
$J_{\text{train}}$  low &  
 $J_{\text{cv}}$  high

$\Rightarrow$  high variance (overfit)

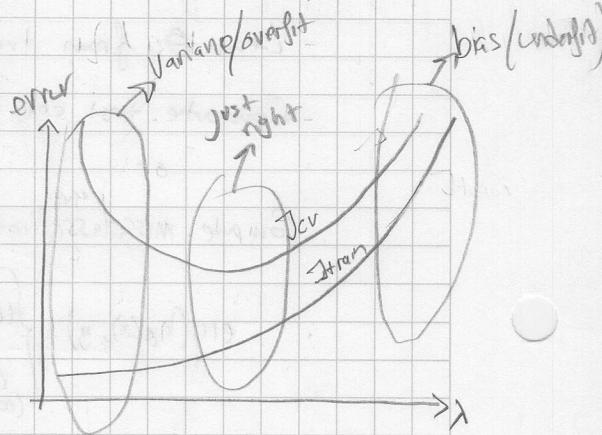
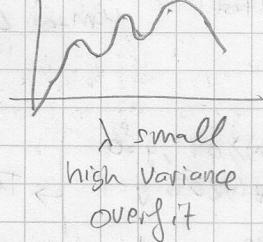
## How to pick the regularization parameter $\lambda$ ?

regularization is used to prevent overfitting...

If  $\lambda$  is too large,  
we 'kill' all  $\theta_1, \dots, \theta_n$   
and  $h_{\theta}(x) = \theta_0$



if  $\lambda$  is too small,  
 $\theta_1, \dots, \theta_n$  are not regularized.



Try:

- (1)  $\lambda = 0 \rightarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$
- (2)  $\lambda = 0.01 \rightarrow \dots \rightarrow \theta^{(2)} \rightarrow J_{\text{cv}}(\theta^{(2)})$
- (3)  $\lambda = 0.02 \rightarrow \dots \rightarrow \theta^{(3)} \rightarrow \dots$
- (4)  $\lambda = 0.04 \rightarrow \dots \rightarrow \theta^{(4)} \rightarrow \dots$
- (5)  $\lambda = 0.08 \rightarrow \dots \rightarrow \theta^{(5)} \rightarrow \dots$
- ⋮
- (12)  $\lambda = 10, 24 \rightarrow \dots \rightarrow \theta^{(12)} \rightarrow J_{\text{cv}}(\theta^{(12)})$

minimum  
 $J_{\text{cv}}(\theta^*) \rightarrow J_{\text{test}}(\theta^*)$  is the error

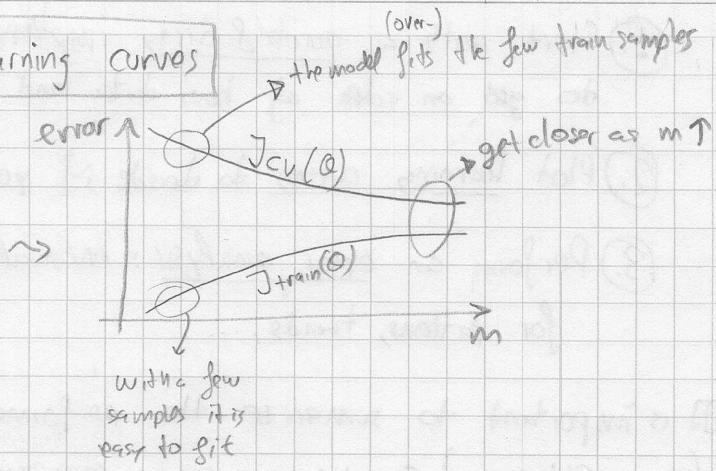
we fit "x" to  
 $J_{\text{cv}}$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum (h_{\theta}(x) - y)^2$$

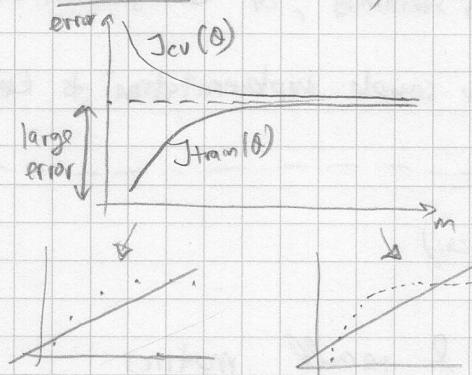
Jcv neglects  $\lambda$  factor (just squared error)

## Diagnosing Bias/Variance with Learning Curves

Δ learning curve plots the error for different sizes of the training set ( $m$ ). That is, using more samples or less samples.

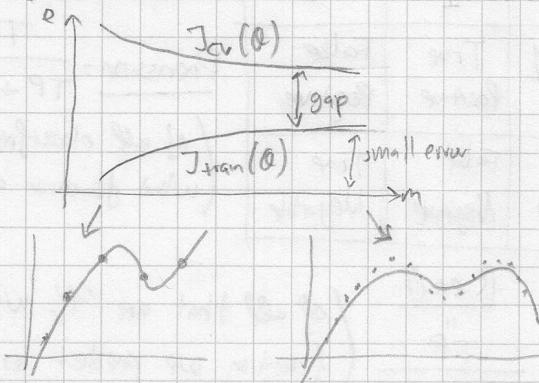


### High bias (underfitting)



} getting more training data will not help that much.  
 $\Rightarrow eq$  does not change with  $m \uparrow$

### High variance (overfitting)



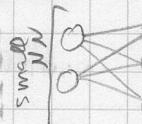
} getting more training data might help  
 $\Rightarrow eq$  closer with  $m \uparrow$

## Revisiting 'What to do next'

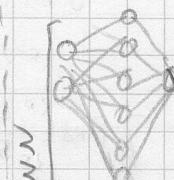
- Get more training data → ✓
- Smaller set of features → ✓
- Get additional features → ✓
- Add polynomial features → ✓
- Decrease  $\lambda$  → ✓
- Increase  $\lambda$  → ✓

high variance (overfit)  
high bias (underfit)

## Back to NN



computationally cheap  
prone to underfitting (high bias)



computationally expensive  
prone to overfitting (high variance)  
use regularization ( $\lambda$ )

How to choose  $L$ ?

- Make train/cv/test splits and test with different architectures

## ML Projects: Recommended Approach

do not spend months on it, in one week you have it.

- ① Start with a quick & dirty implementation of a simple algorithm to get an idea of how data and learning behaves (in CV)
- ② Plot learning curves to decide if you need more data, more features, etc...
- ③ Perform an error analysis: manually examine missclassified examples looking for patterns, trends, ...

It is important to summarize the performance into a single numerical value (e.g. CV error) so we can easily compare different strategies. For instance, in a spam classifier, compare 'use stemming' vs 'not using stemming', or 'distinguish Capital!' vs 'not distinguishing'.

→ try different ideas with the simple implementation to keep improving the algorithm.

## Alternative Error Metrics

(good for skewed classes)

When we have a classifier (with 2 classes) where in the data there is only 1% of one of the type classes, it is difficult to measure learning accuracy with common error rates, because it is easy to reach ~99% of correct classifications but just due to data skewness.

Use precision & recall metrics

		Actual class		$\frac{\text{true positives}}{\text{predicted positive}}$
		1	0	
class	1	True Positive	False Positive	Precision = $\frac{\text{TP}}{\text{TP} + \text{FP}}$ (of all classified as '1', what fraction actually was '1')
	0	False Negative	True Negative	

$$\text{Recall} = \frac{\text{true pos.}}{\text{actual positive}} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{of all that are '1', what fraction we predict as '1'})$$

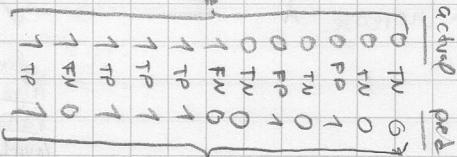
We want high precision &

high recall



prediction 1 0 0 0 0 0 1 1 1 1 1      Pred =  $\frac{4}{5} = 0.8$   
 actual 0 0 1 0 1 1 0 1 1 1      Recall =  $\frac{4}{5} = 0.66$

Recall → of all that were '1', how many were predicted '1'



of all predicted '1', how many were really '1'

← Precision

→ but usually there is a tradeoff between precision and recall.

For example, in a log Reg we can set a <sup>high</sup> ~~low~~ threshold:

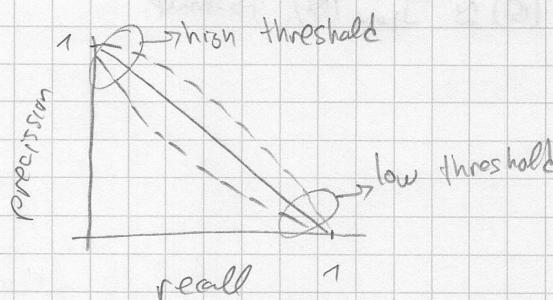
$y=1$  if  $h_0(x) > 0.8$  → we are more confident of the  $y=1$  predictions

(A)  $y=0$  if  $h_0(x) < 0.8$  → higher precision, lower recall  
$$\left( \frac{TP \downarrow}{\text{predicted Pos} \uparrow} \right) \quad \left( \frac{TP \downarrow}{\text{actual pos.} = \uparrow} \right)$$

or we can set a <sup>low</sup> ~~high~~ threshold:

$y=1$  if  $h_0(x) > 0.3$  → we avoid missing positive classifications ( $\downarrow FN$ )

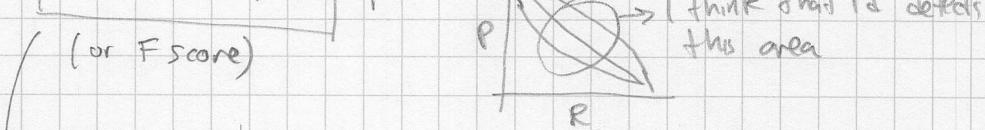
(B)  $y=0$  if  $h_0(x) < 0.3$  → lower precision, higher recall  
$$\left( \frac{TP \uparrow}{\text{predicted Pos} \uparrow} \right) \quad \left( \frac{TP \uparrow}{\text{actual pos.} = \uparrow} \right)$$



The problem with using Precision & Recall (including the tradeoff) is that we have lost the 'single metric' evaluation that we had with the CV error. How can we transform P&R into a single metric?

~~Average =  $\frac{P+R}{2}$~~  → bad metric because 'bad' classifiers with very high R & very low P (or vice versa) end up with medium average...

$$F_1 \text{ score} = 2 \frac{PR}{P+R} \rightarrow \text{better measures when P & R are large...}$$



Use this in the CU set to choose among different algorithms.

## Large Data Rationale

The key question is when it makes sense to collect a large amount of data. This is an important issue because not in all conditions having more data results in more accuracy.

The rationale says that:

1- First use a learning algorithm with many parameters

→ low bias, prone to high variance

→  $J_{\text{train}}(\theta)$  will be small

2- Then reduce the variance with more data

→ low bias, low variance

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