

①

Machine Learning - Week 6 - Advice for Applying Machine Learning

T1 - Evaluating a Learning Algorithm

L1 - Deciding what to try next

Debugging a learning algorithm:

Suppose you have implemented a regularized linear regression to predict housing prices.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding Polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc...)
- Try decreasing λ
- Try increasing λ

Machine Learning diagnosis:

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

Diagnostics can take time to implement, but doing so can be a very good use of your time.

L2 - Evaluating a hypothesis

Dataset:	Size	Price	
	2104	400	
	1600	330	
	2400	368	
70%	1416	232	} Training set →
	3000	540	
	1985	300	
	1534	315	
	1427	199	
	1380	212	} Test set →
30%	1494	243	

Select Random examples

$(x^{(1)}, y^{(1)})$

$(x^{(2)}, y^{(2)})$

...

$(x^{(m)}, y^{(m)})$

$(x_{test}^{(1)}, y_{test}^{(1)})$

$(x_{test}^{(2)}, y_{test}^{(2)})$

...

$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$M_{test} = \# \text{ of test examples}$

Training/testing procedure for linear regression

- Learn parameter θ from training data (minimize training error $J(\theta)$)

- Compute test set error:

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

Training/testing procedure for logistic regression

- Learn parameter θ from training data

- Compute test set error:

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

- Misclassification error (0/1 misclassification error):

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

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

13- Model Selection and training/validation/test sets

$d = \text{degree of polynomial}$

Model Selection

$$\begin{array}{llll} d=1 & 1. & h_{\theta}(x) = \theta_0 + \theta_1 x & \xrightarrow{\theta^{(1)}} J_{\text{test}}(\theta^{(1)}) \\ d=2 & 2. & h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 & \xrightarrow{\theta^{(2)}} J_{\text{test}}(\theta^{(2)}) \\ d=3 & 3. & h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 & \xrightarrow{\theta^{(3)}} J_{\text{test}}(\theta^{(3)}) \\ \vdots & \vdots & \vdots & \vdots \\ d=10 & 10. & h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} & \xrightarrow{\theta^{(10)}} J_{\text{test}}(\theta^{(10)}) \end{array}$$

Select min

Ex. choose $\theta_0 + \dots + \theta_5 x^5$ ($d=5$)

How well does the model generalize? Report test set error $J_{\text{test}}(\theta^{(5)})$.

Problem: $J_{\text{test}}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error, i.e. our extra parameter ($d = \text{degree of polynomial}$) is fit to test set.

②

Evaluating your hypothesis

Dataset:

	Size	Price		
60%	2104	400	Training set	$(x^{(1)}, y^{(1)})$ \vdots $(x^{(m)}, y^{(m)})$
	1600	330		
	2400	369		
	1416	232		
	3000	540		
20%	1975	300	Cross-validation set (CV)	$(x_{cv}^{(1)}, y_{cv}^{(1)})$ \vdots $(x_{cv}^{(m)}, y_{cv}^{(m)})$
	1534	315		
	1427	199		
20%	1380	212	Test set	$(x_{test}^{(1)}, y_{test}^{(1)})$ \vdots $(x_{test}^{(m)}, y_{test}^{(m)})$
	1494	243		

Train/Validation/Test error.

Training error:

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

Cross Validation error:

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

Test Error:

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

Model Selection

$$1. h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^* \rightarrow J_{\text{cv}}(\theta^{(1)})$$

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

$$\text{e.g. Pic 11 } \theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4 \quad (d=4)$$

Estimate generalization error for test set $J_{\text{test}}(\theta^{(4)})$

#2- Bias vs Variance

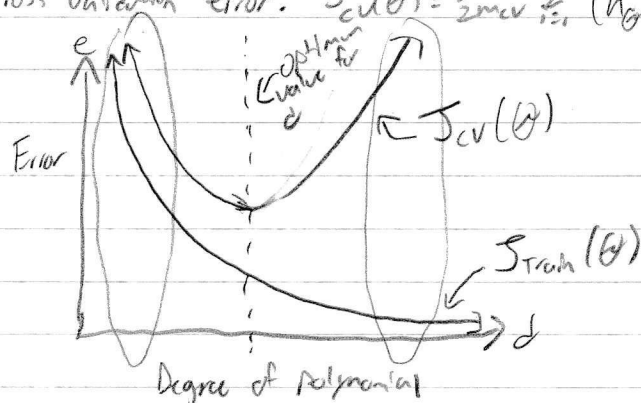
L1- Diagnosing Bias vs Variance

Recall:

High bias (underfit) and High variance (overfit)

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

$$\text{Cross Validation error: } J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$



Suppose your learning algorithm is performing less well than you were hoping.
($J_{\text{cv}}(\theta)$ or $J_{\text{test}}(\theta)$ is high). Is it a bias problem or a variance problem?

Bias (underfit)

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

Variance (overfit)

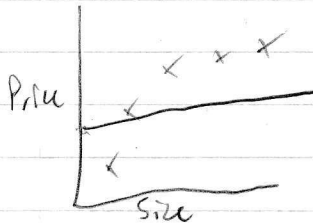
- $J_{\text{train}}(\theta)$ will be low
- $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$

L2- Regularization and bias/variance

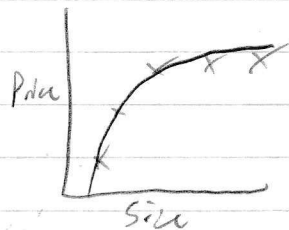
Linear regression with regularization

$$\text{Model: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

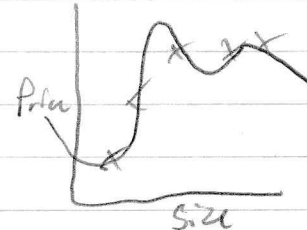
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^4 \theta_j^2$$



Large λ
High bias (underfit)
 $\lambda = 10,000, \theta_1 \approx 0, \theta_2 \approx 0, \dots$
 $h_{\theta}(x) \approx \theta_0$



Intermediate λ
"just right"



Small λ
High variance (overfit)
 $\lambda = 0$

Choosing the regularization parameter λ

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

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

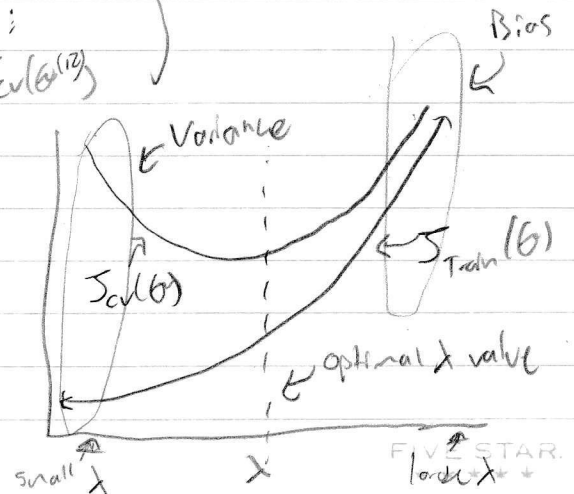
$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

Note there is no λ

1. Try $\lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$
2. Try $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{\text{cv}}(\theta^{(2)})$
3. Try $\lambda = 0.02 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(3)} \rightarrow J_{\text{cv}}(\theta^{(3)})$
4. Try $\lambda = 0.04 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(4)} \rightarrow J_{\text{cv}}(\theta^{(4)})$
- \vdots
12. Try $\lambda = 10 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(12)} \rightarrow J_{\text{cv}}(\theta^{(12)})$

Select min

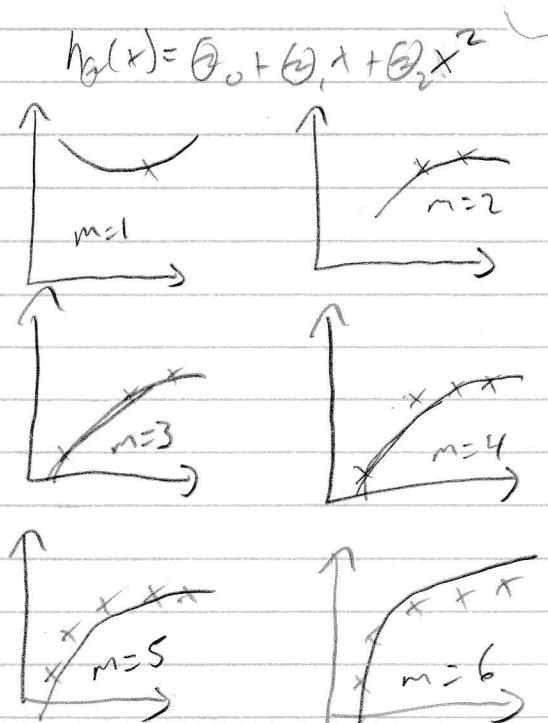
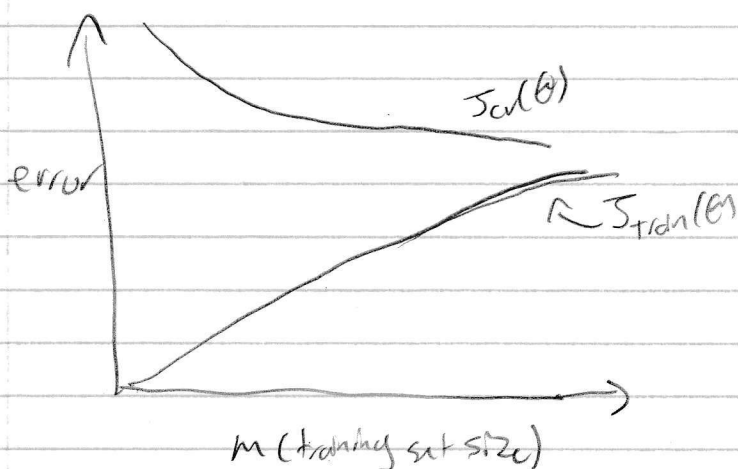
eg. Pick $\theta^{(5)}$ Test error: $J_{\text{test}}(\theta^{(5)})$



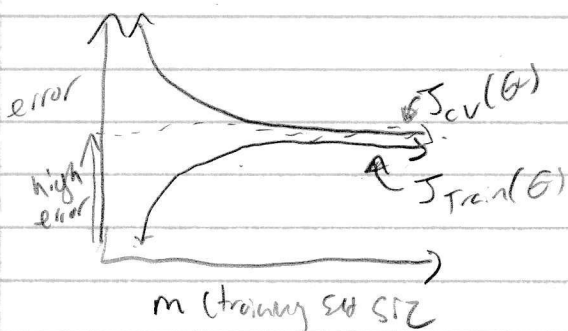
L3- Learning Curves

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

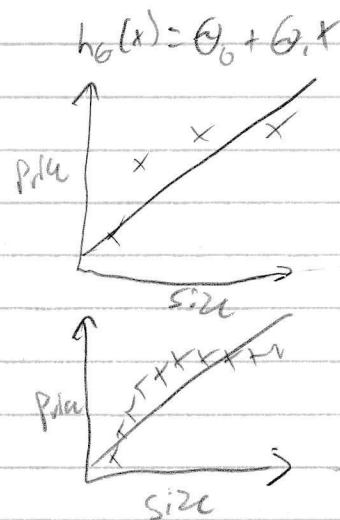
$$J_{\text{dev}}(\theta) = \frac{1}{2m_{\text{dev}}} \sum_{i=1}^{m_{\text{dev}}} (h_{\theta}(x_{\text{dev}}^{(i)}) - y_{\text{dev}}^{(i)})^2$$



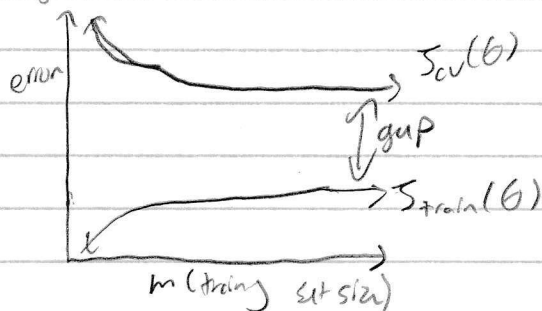
High Bias



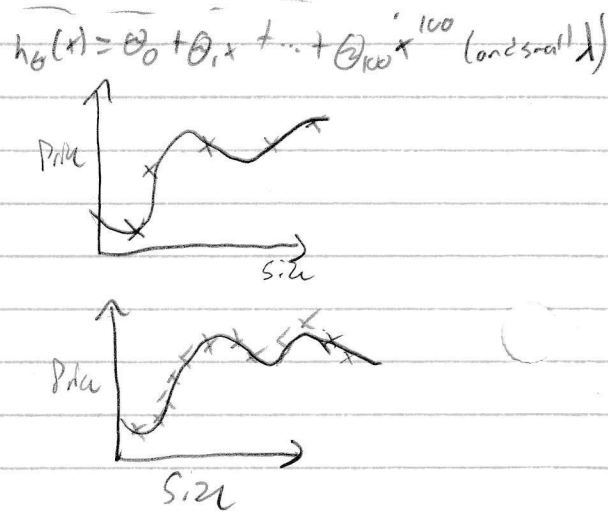
If a learning algorithm is suffering from high bias, getting more data will not (by itself) help much.



High Variance



If a learning algorithm is suffering from high variance, getting more training data is likely to help.



L4- Deciding what to try next (revisited)

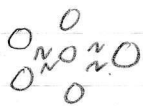
Recall:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Get more training examples
 - Fixes high variance
- Try smaller sets of features
 - Fixes high variance
- Try to get additional features
 - Fixes high bias
- Try adding polynomial features (x_1^2, x_2^2, x_1x_2 , etc)
 - Fixes high bias
- Try decreasing λ
 - Fixes high bias
- Try increasing λ
 - Fixes high variance

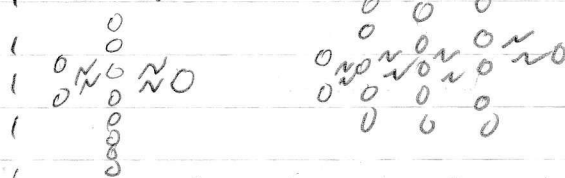
Neural Networks and Overfitting

"Small" nn (few parameters; more prone to underfitting)



- Computationally cheaper

"Large" nn (more parameters; more prone to overfitting)



- Computationally more expensive
- Use regularization (λ) to address overfitting

Machine Learning - Week 6 - Machine Learning System Design

T3 - Building a Spam Classifier

L1 - Prioritizing what to work on

System Design Example:

Given a data set of emails, we could construct a vector for each email. Each entry in this vector represents a word. The vector normally contains 10,000 to 50,000 entries gathered by finding the most frequently used words in our data set. If a word is to be found in the email, we would assign its respective entry a 1, else that entry would be a 0. Once we have all our x vectors ready, we train our algorithm and finally, we could use it to classify if an email is spam or not.

Building a Spam Classifier.

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

Feature X : Choose 100 words indicative of spam/not spam

So how could you spend your time to improve the accuracy of the classifier?

- Collect lots of data (e.g. honey pot project but doesn't always work)
- Develop sophisticated features (e.g. using email header data in spam emails)
- Develop algorithms to process your input in different ways (recognizing misspellings in spam)

It is difficult to tell which of the options will be most helpful.

L2 - 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 examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of example it is making errors on.

T4 - Handling Skewed Data

L1 - Error metrics for Skewed Classes

Cancer class. Another example

Train logistic regression model $h_0(x)$. ($y=1$ if cancer, $y=0$ otherwise)

Find that you get 1% error on test set. (99% correct diagnosis)

But only 0.50% of patients have cancer.

↳ skewed classes.

function $y = \text{predictCancer}(x)$ } 0.5 % error
 $y=0$; % ignore x !
 return

Precision/Recall

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

Actual class

Predicted class	1	0
	True Positive False Positive	False Positive True Negative
0	False Negative True Negative	True Negative True Negative

Precision (of all the patients where we predicted

$y=1$, what fraction actually has cancer?

i.e. True positives

True positives + False positives

↳ Predicted positives

High Precision and high Recall is good

Recall (of all patients that actually have cancer, what

fraction did we correctly detect as having cancer?

i.e. True positives

True pos + False neg

↳ Actual positives

L2 - Trading off Precision and Recall

Logistic regression: $0.5 \leq h_0(x) \leq 1$

Predict 1 if $h_0(x) \geq 0.5$

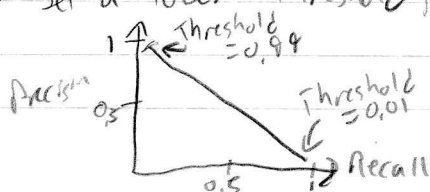
Predict 0 if $h_0(x) < 0.5$

Suppose we want to predict $y=1$ (cancer) only if very confident.

- Could set a higher threshold, resulting in higher precision and lower recall

Suppose we want to avoid missing too many cases of cancer (avoid false negatives).

- Could set a lower threshold, resulting in higher recall and lower precision



F₁ Score (F score)

How to compare precision/recall numbers?

	Precision (P)	Recall (R)	Average	F ₁ score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.51	0.0392

Average: $\frac{P+R}{2}$
- Not very effective

F₁ score: $2 \frac{P \cdot R}{P+R}$

$P=0$ or $R=0 \Rightarrow F\text{-score} = 0$

$P=1$ and $R=1 \Rightarrow F\text{-score} = 1$

IS - Using Large Data Sets

LI - Data for Machine Learning

Large data failure

Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately

Ex. For breakfast I ate _____ eggs. (two, too, or fo)?

Counter Ex. Predict housing price from only size (feet²) and no other features.

Useful test: Given the input x , can a human expert confidently predict y ?

Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features, neural network with many hidden units). Low bias algorithms.

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

Use a very large training set (unlikely to overfit). Low variance

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

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