

Debug and improve learning algorithm.

some promising avenue :

- | | |
|-----------------------------------|---------------------|
| Get more training examples | → fix high variance |
| trying a smaller sets of features | → fix high variance |
| trying additional features | → fix high bias |
| adding polynomial features | → fix high bias |
| increasing λ | → fix high variance |
| decreasing λ | → fix high bias |

a hypothesis may have low error for the training examples but still be inaccurate, because of overfitting, thus we split the dataset into two sets : a training set and a test set

training set (70%) test set (30%)
learn θ and minimize $J_{\text{train}}(\theta)$ → compute the test set error $J_{\text{test}}(\theta)$

The Test Set Error :

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

for classification : misclassification error aka 0/1 misclassification error

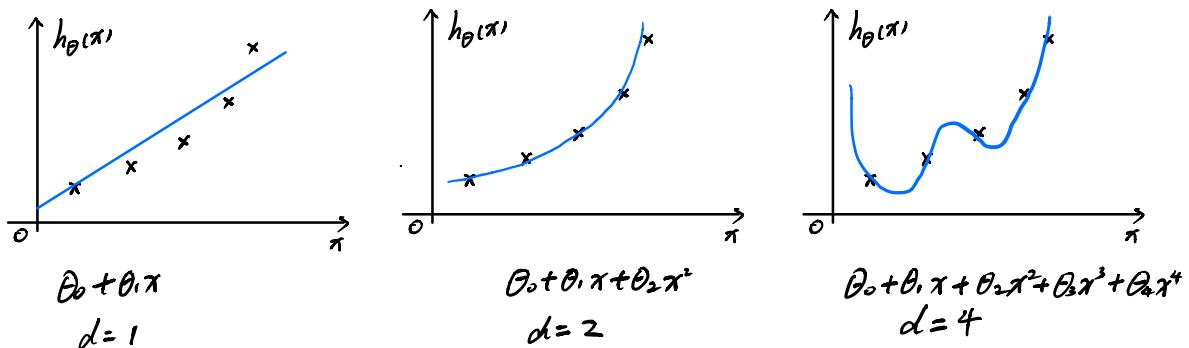
$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) > 0.5 \text{ and } y = 0 \text{ or } h_{\theta}(x) \leq 0.5 \text{ and } y = 1 \\ 0 & \text{otherwise} \end{cases}$$

the average test error is :

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

which gives us the proportion of the test data be misclassified.

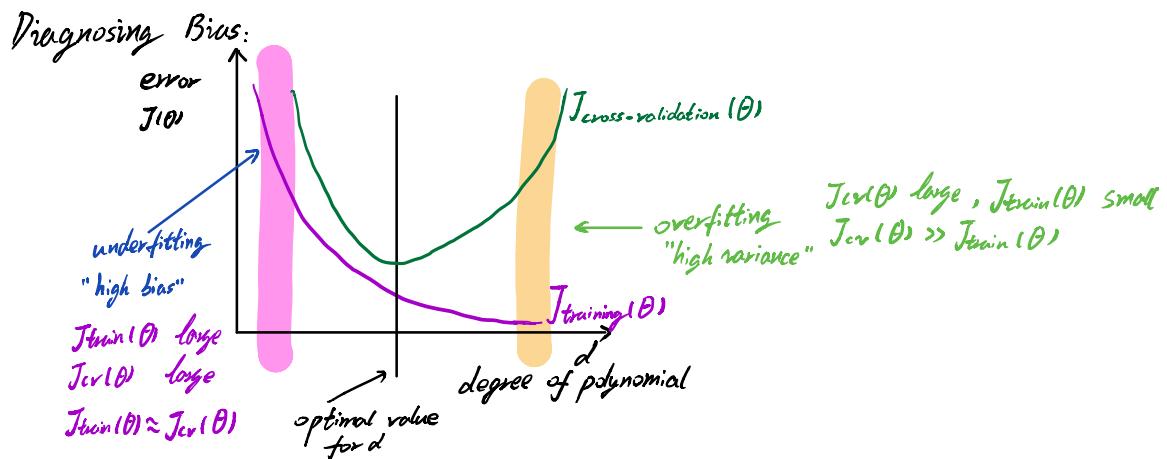
Model Selection and train/validation/test sets
 choose the degree of polynomial function to decide the model.



60% <u>training Set</u>	20% <u>Cross validation set</u>	20% <u>test set</u>
-------------------------------	---------------------------------------	---------------------------

learn θ and minimize $J(\theta)$ for each polynomial degree, \Rightarrow least error using cross validation set. $\min_{\theta} J_{cv}(\theta) = J(\theta^{cv})$

select the degree d with estimate the generalization error using the test set with $J_{test}(\theta^{cv})$

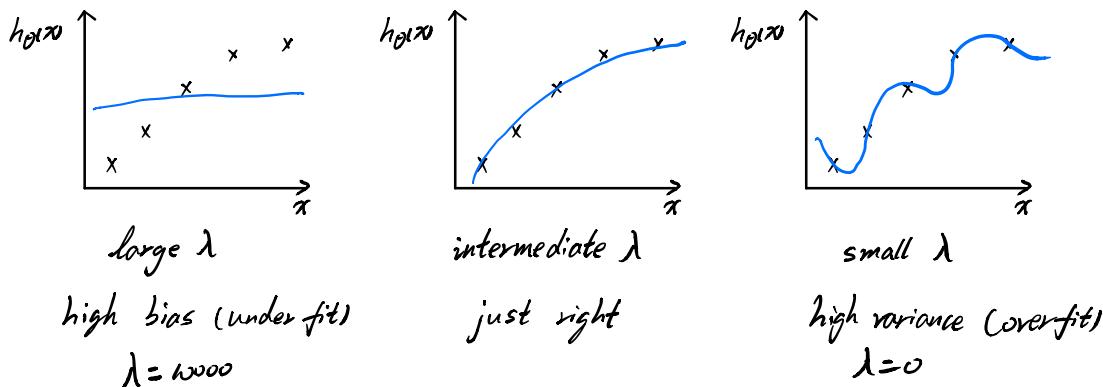


Regularization and Bias/Variance

choosing the regularization parameter λ

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

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



split the dataset by 60%, 20%, 20%, compute $\min_k J_{cv}(\theta^{(k)})$, then select λ

try $\lambda = 0$	x_2	\rightarrow	fit $\min_{\theta} J(\theta) \rightarrow \theta^{(0)}$	\rightarrow	compute $J_{cv}(\theta^{(0)})$	}
$\lambda = 0.01$	x_2	\rightarrow	$\min_{\theta} J(\theta) \rightarrow \theta^{(1)}$	\rightarrow	$J_{cv}(\theta^{(1)})$	
$\lambda = 0.02$	x_2	\rightarrow	$\min_{\theta} J(\theta) \rightarrow \theta^{(2)}$	\rightarrow	$J_{cv}(\theta^{(2)})$	
$\lambda = 0.04$	x_2	\rightarrow	$\min_{\theta} J(\theta) \rightarrow \theta^{(3)}$	\rightarrow	$J_{cv}(\theta^{(3)})$	
\vdots		\vdots		\vdots		
$\lambda = 10$		\rightarrow	$\min_{\theta} J(\theta) \rightarrow \theta^{(10)}$	\rightarrow	$J_{cv}(\theta^{(10)})$	

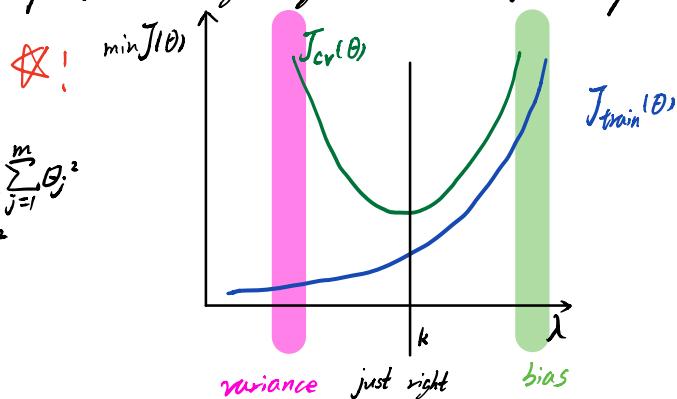
then compute $J_{test}(\theta)$ to see if it has a good generalization of the problem.

if let:

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

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

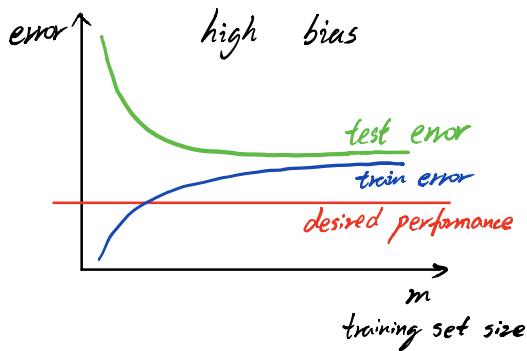
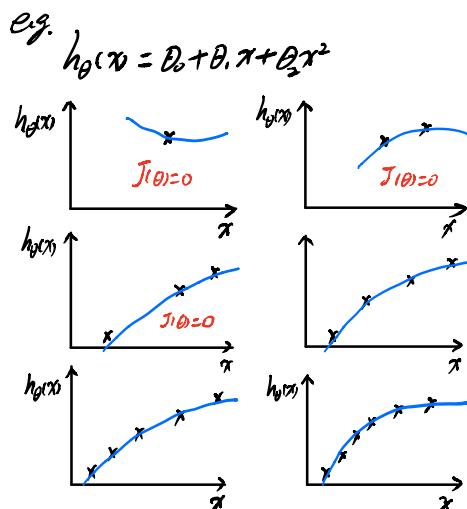
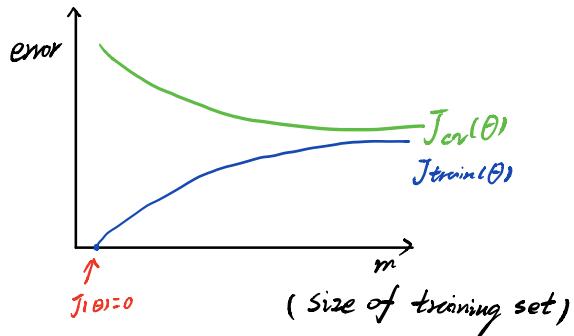
$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



Learning Curves

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

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



m large, $J_{\text{train}}(\theta) \approx J_{\text{cv}}(\theta)$

getting more training data doesn't help

m large, $J_{\text{train}}(\theta) < J_{\text{cv}}(\theta)$

getting more training data helps

Diagnose neural networks

neural network with fewer parameters, layers is prone to underfitting, it's also computationally cheap.

neural network with more parameters, layers is prone to overfitting, it's also computationally expensive: you can use regularization to address fitting.

neural network with multiple hidden layer using regularization prone to perform better than single hidden layer network. You can train multiple hidden layer network using cross validation set and then choose the one perform best.

Machine Learning System design.

Build a spam classifier

a supervised learning example, x = features of the email, y = spam(1) or not spam(0).

feature x : choose n ($n > 100$) words indicative of spam/not spam.

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{andrew} \\ \text{buy} \\ \text{deal} \\ \text{discount} \\ \vdots \\ \text{now} \end{array} \quad x \in \mathbb{R}^n \quad x_i = 1 \text{ if } i\text{th word appears}$$

how to spend your time to make it have low error:

- collect lots of data
- develop sophisticated features based on email routing information or message body
- develop sophisticated algorithm to detect misspelling.

Recommended approach on starting a machine learning problem:

- start with a simple algorithm that you can implement quickly. Implement it and test it on your cross-validation dataset.
- plot learning curves to decide more data, more features etc. may help.
- error analysis: manually examine the examples in cross-validation set that your algorithms made errors on. See if you spot any systematic trend in what type of example it is making errors on.

Error Analysis :

$M_{cv} = 500$, algorithm misclassified 100 emails, manually examine the 100 errors, and categorize them based on:

(i) what types of the misclassified email? $\Rightarrow \begin{cases} \text{pharma: } 12 \\ \text{replica/fake: } 4 \\ \text{steal password: } 53 \\ \text{other: } 31 \end{cases}$

(ii) what features you think would have helped algorithm classify them correctly? $\Rightarrow \begin{cases} \text{misspelling: } 5 \\ \text{unusual email routing: } 16 \\ \text{unusual punctuation: } 32 \end{cases}$

need numerical evaluation of algorithm's performance:

distinguish upper/lower case? \checkmark 3% error N 5% error

Error Metrics for skewed classes:

e.g. cancer classification problem:

use a logistic regression model $h_\theta(x)$ to predict cancer (1), otherwise (0)
find that you got 1% error on test set (99% correctly diagnoses)
actually, only 0.5% of patients have cancer.

if we define function predicting cancer: $y=0$, means for all cases just return 0 (not cancer), the accuracy of hypothesis should be 99.5%, and the error rate = $1 - 99.5\% = 0.5\%$, it looks pretty good but it is obviously not a proper hypothesis function. Accuracy should not be used to measure the performance of hypothesis for they are skew classes (the difference of cancer and not cancer is huge), hence we need to introduce a new measure:

Precision / Recall

we usually set $y=1$ in presence of rare class (cancer)

		actual class	
		1	0
predicted class	1	true positive	false positive
	0	false negative	true negative

Precision

(of all patients where we predicted $y=1$, the fraction actually has cancer)

$$= \frac{\text{true positive}}{\text{predicted positive}} = \frac{\text{true positive}}{\text{true positive} + \text{false positive}}$$

Recall

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

$$= \frac{\text{true positive}}{\text{actually positive}} = \frac{\text{true positive}}{\text{true positive} + \text{false negative}}$$

$$\text{Accuracy} = \frac{\text{true positive} + \text{true negative}}{\text{all examples}}$$

Trade off precision and recall

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

- predict 1 if $h_0(x) > \text{threshold}$ (usually 0.5)
- predict 0 if $h_0(x) < \text{threshold}$

to select the most suitable threshold, we need to trade off between precision and recall, and use a single number metric.

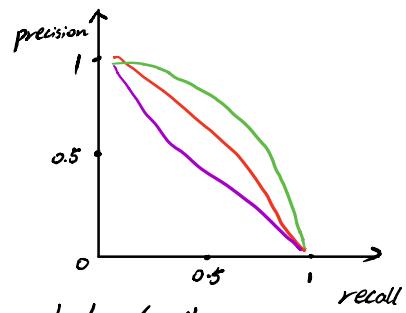
Suppose we want to predict $y=1$ (cancer) only if very confident:

threshold = 0.7 ↑ precision ↑ recall ↓

Suppose we want to avoid missing too many cancer:

threshold = 0.3 ↓ precision ↓ recall ↑

e.g.	Precision (P)	Recall (R)
Algorithm 1	0.5	0.4
Algorithm 2	0.7	0.1
Algorithm 3	0.02	1.0



if y always = 1 actually this is a bad algorithm,
so we can't use average $\frac{R+P}{2}$ to measure.

we use F_1 score: $2 \frac{PR}{P+R}$

when $P=0$ or $R=0$, F_1 score = 0

$P=1$ and $R=1$, F_1 score = 1

Data For Machine Learning

— how many data to use.

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

I: useful test: given the input test x , can a human expert confidently predict y ?

II: use a learning algorithm with many parameters (many features, many hidden layers)

→ this tends to be a low bias algorithm

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

III. use a large training set (unlikely to overfit)

→ this tends to be a low variance algorithm

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

} meet both can get a good algorithm