

## Evaluating a Learning Algorithm

Deciding what to try next

### Decisions

find it makes unacceptably large errors.  
What should you try next?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features  
 $(x_1^2, x_2^2, x_1 x_2, \text{etc})$
- Try decreasing / increasing  $\lambda$

### Definition of ML 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.

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

## Evaluating a hypothesis

Overfitting judgement. (train error 小)  
(test error 大)

將 Dataset 上以 7:3 的比例隨機拆成  
兩大部分，70% 作為訓練集  
30% 作為測試集

### For linear regression

- Learn parameter  $\theta$  from training data (minimizing training error  $J(\theta)$ )
- Compute test set error :

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

## For logistic regression

- Learn parameter
- Compute test set error :

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

- Misclassification error (%, misclassification error)

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

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

→ This gives us the proportion of the test data that was misclassified.

## Model Selection and Train / Validation / Test Sets

### Model Selection

以訓練集訓練

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

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

⋮  
⋮  
⋮

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

比較  $\theta^{(1)} \sim \theta^{(10)}$  的 test error 值之大小  
選擇最小的 error 值之  $\theta$  當作  $h_{\theta}$

聽起來很合理？

但，此  $\theta$  是看過測試集的，因此，不能代表此  $\theta$  對新樣本的擬合能力。  
也比其他  $\theta$  表現的更加優異  
*Selecting your model using this test set,  
and then using the same test set to report  
the error as though!?*

那，該怎麼辦呢？

將 Dataset 分成訓練集、測試集、驗証集  
*Cross validation*

比例為 60% : 20% : 20%

解決方法：用驗証集進行目的挑選而不是用測試集

### Train / Validation / Test Sets

Training error:

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

用途：用來訓練模型。確定模型參數

Cross Validation error:

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

用途：用來做模型選擇

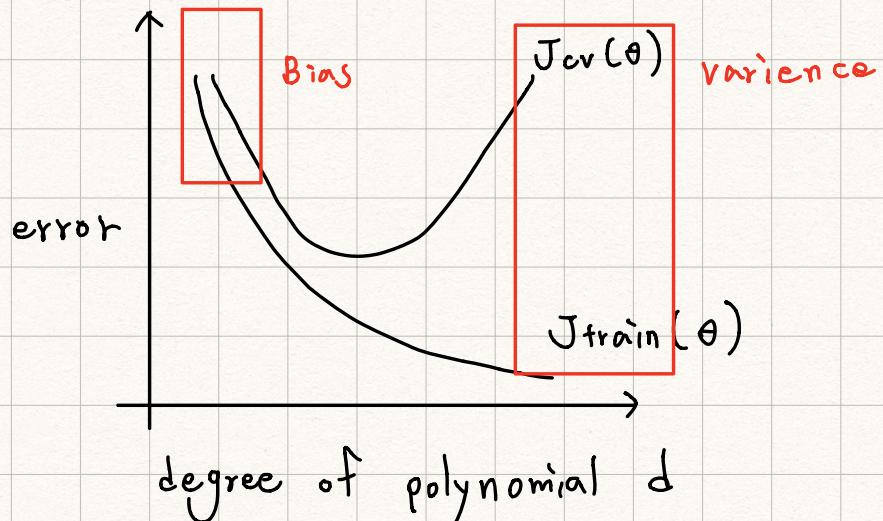
Test error:

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

用途：純粹為了測試已經訓練好的模型之泛化能力

## Diagnosing bias v.s. variance

評價算法存在 bias or variance 的問題



Bias :

$J_{cv}(\theta)$  及  $J_{train}(\theta)$   
都偏大且數值接近  
Underfit

Variance :

$J_{cv}(\theta)$  偏大,  $J_{train}(\theta)$  偏小  
表其 overfit

## Regularization and Bias / Variance

model:  $h(\theta) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$$J(\theta) = \frac{1}{m} \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

$$\left\{ \begin{array}{l} J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ J_{\text{cv}}(\theta) = \frac{1}{m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2 \\ J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2 \end{array} \right\}$$

Choosing the regularization parameter  $\lambda$

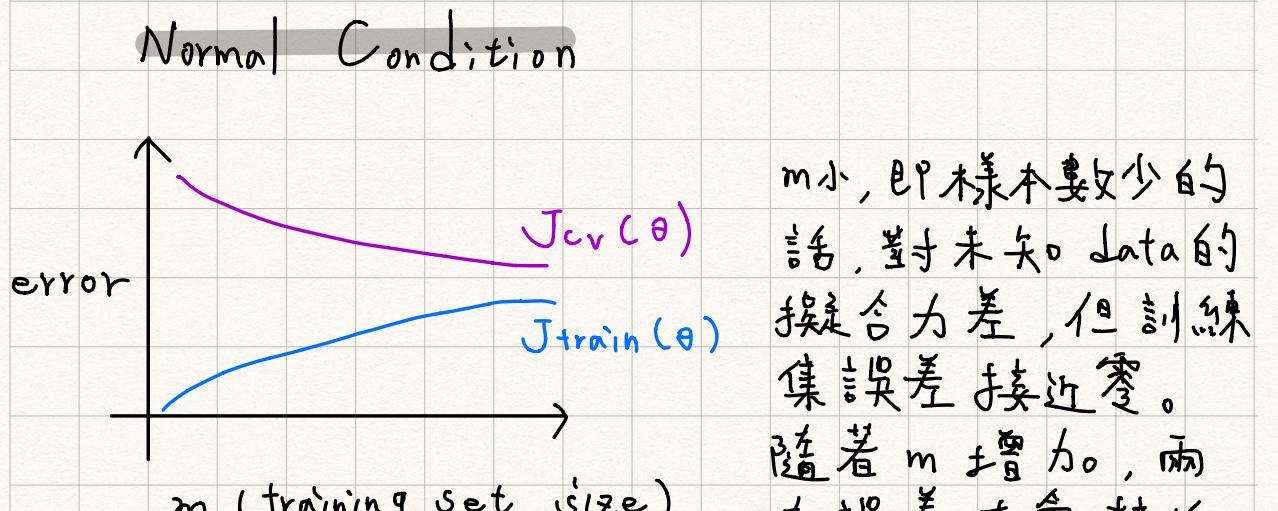
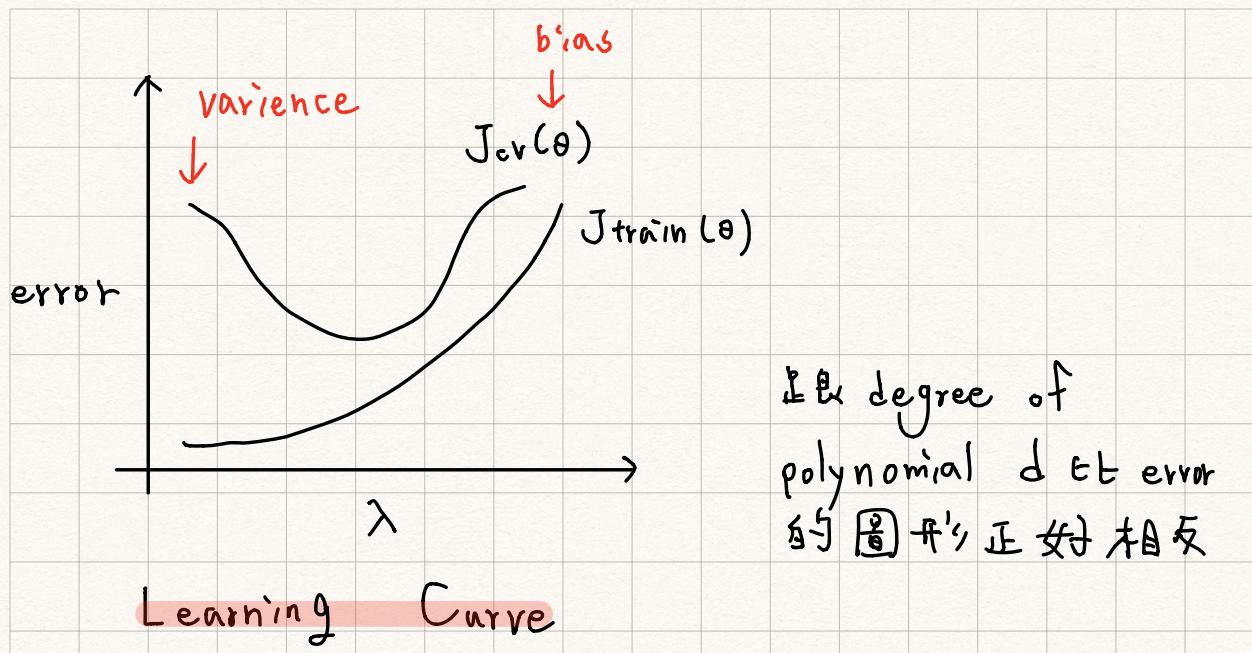
|         |                  |  |                                |
|---------|------------------|--|--------------------------------|
| 1. Try  | $\lambda = 0$    | $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)}$ | $J_{\text{cv}}(\theta^{(1)})$  |
| 2. Try  | $\lambda = 0.01$ | $\dots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

Pick  $\theta$  so that Test Error is the lowest ( $\lambda$ )

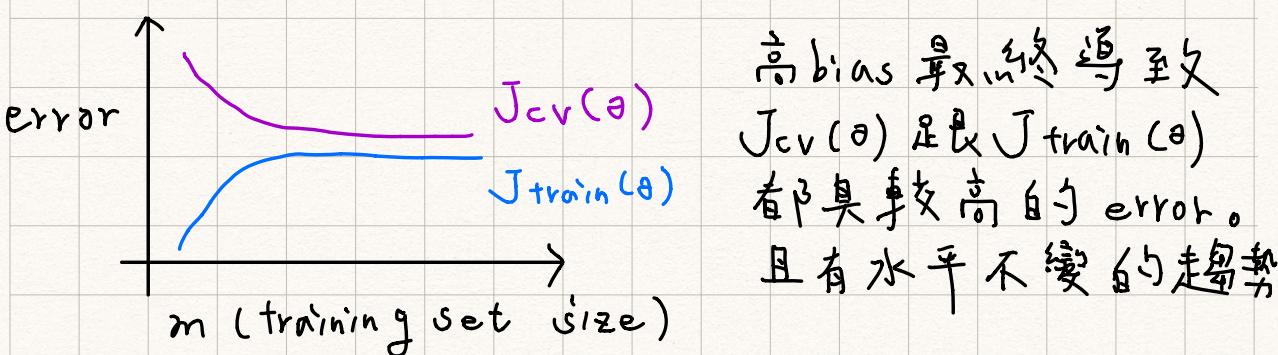
$\lambda \rightarrow$  underfit,  $J$  small  $\rightarrow$  overfit,  $J$  large  $\rightarrow$  perfect

Bias/variance as a function of the regularization

$x: \lambda$ ,  $y: \text{error}$

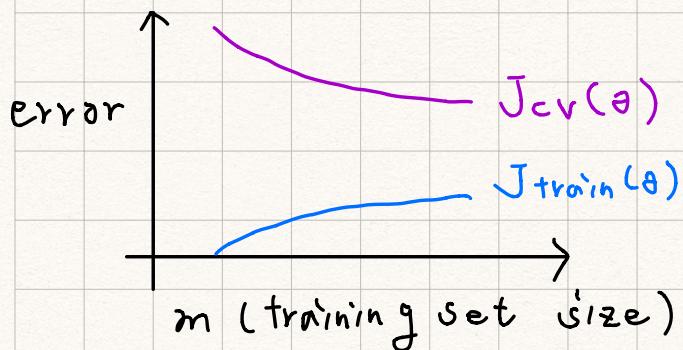


High bias



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

High variance



高 variance 導致  
在相同的 m 下，  
 $J_{cv}(\theta)$  和  $J_{train}(\theta)$   
的差比其他狀況  
的差值大得多

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

### Deciding What to Do Next Revisited

find it makes unacceptably large errors.  
what should you try next?

- 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  
 $(x_1^2, x_2^2, x_1x_2, \text{etc})$   
→ fixes high bias
- Try decreasing / increasing  $\lambda$   
→ fixes high variance

## Neural networks and overfitting

### "Small" neural network

- fewer parameters ; more prone to underfitting
- Computationally cheaper

### "Large" neural network

- more parameters ; more prone to overfitting
- Computationally more expensive

# Machine Learning System Design

## Building a spam classifier

### Prioritizing What to Work on

找出該系統的特徵 (x or y)

Choose 100 words indicative of spam/not spam

$$x = \begin{bmatrix} 0 \\ 1 \\ -1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \begin{array}{l} \text{andrew} \\ \text{buy} \\ \text{deal} \\ \text{discount} \\ \vdots \\ \text{now} \end{array}$$

→ take most frequently occurring n words (10000 in training set, ~50000) rather than manually pick 100 words.

How to spend your time to make it have low error

- Collect lots of data
- Develop sophisticated features based on email routing information
- Develop sophisticated features for message body
- Develop sophisticated algorithm to detect misspellings

## Error Analysis

### Recommended approach

- Start with a simple algorithm
- Plot learning curve
- Error analysis

### Error Analysis Example

500筆馬郵証集中，有100筆子員測錯誤  
解決方法？

(i) What type of email it is.

pharma: 12      Replica / Fake : 4  
→ steal passwords : 53      Other: 31

(ii) What cues (features) you think would have helped the algorithm classify them correctly.

Deliberate misspelling : 5  
Unusual email routing : 16  
→ Unusual punctuation : 32

人工搜尋了臭測錯誤的 100 節 mail 中，  
屬於哪種類型的最多

或是，那些特徵在信件中是最常出現  
的（即對臭測最有幫助的）

The importance of numerical evaluation  
Error analysis may not be helpful for  
deciding if this is likely to improve  
performance. Only solution is to **try**  
**it and see** if it works.

Note: 在驗證集上做誤差分析

一開始不必以複雜的算法進行  
訓練，而以簡單快速的算法去  
干。等訓練完成後，做誤差分  
析看此算法出了什麼錯，然  
後再決定往後優化的方向。

再加上適當的數值評估依據，  
便能幫助你嘗試新方法，以及  
這些方法是否有交叉用

## Error metrics for skewed classes

Cancer classification example

① logistic regression model.

Find that you got 1% error on test set.

Only 0.5% of patients have cancer.

己. 總是預測病人沒得到 cancer

0.5% error

比較起來，己還比己的準確度高

此時，這個 0.5% 便稱為 skewed classes

Precision / Recall

Actual class

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

Precision: Of all patients where we predicted  $y=1$ , what fraction actually has cancer?

Recall : Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?

Precision :

$$\frac{\text{True positives}}{\# \text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

Recall :

$$\frac{\text{True positives}}{\# \text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

Precision 與 Recall 可以用來評估  
算法是否優良

(Precision 與 Recall 都高，表其很 OK。  
不會有總是預測  $y=1$  v  $y=0$  的算法，  
Precision - Recall 都高的情形發生。)

Trading off precision and recall

Logistic regression :  $0 \leq h_\theta(x) \leq 1$

1 : if  $h_\theta(x) \geq 0.5$ , 0 : if  $h_\theta(x) < 0.5$

0.9

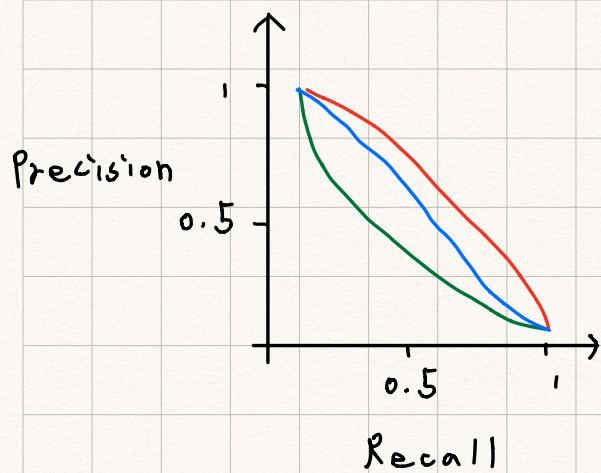
0.1

Higher precision, lower recall

0.3

0.3

Higher recall, lower precision



要高準確 / 召回率，端看  
使用者要怎麼做

More generally:  
Predict 1 if  
 $h_\theta(x) \geq \text{threshold}$ .

## F<sub>1</sub> Score (F score)

|    | Precision (P) | Recall (R) | Average | F <sub>1</sub> Score |
|----|---------------|------------|---------|----------------------|
| A1 | 0.5           | 0.4        | 0.45    | 0.444                |
| A2 | 0.7           | 0.1        | 0.4     | 0.175                |
| A3 | 0.02          | 1.0        | 0.51    | 0.0392               |

Predict  $y=1$  all the time

$$F_1 \text{ Score} = 2 \cdot \frac{P \cdot R}{P + R}$$

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

$P = 1 \text{ or } R = 1 \Rightarrow F\text{Score} = 1$

選擇不同的 threshold,

並在驗證集上測試  $F_1$  的值，以得到品質較好的預測。

## Data for Machine Learning

Designing a high accuracy learning system

2001 年，Banko and Brill 估計的研究指出  
越多數據，預測的結果越好

→ It's not who has the best algorithm  
that wins. It's who has the most data.

Large data rationale

1. Use a Learning algorithm with many  
parameters (low bias algorithms → 拟合效果好)

■  $J_{train}(g)$  will be small (low bias)

2. Use a very large training set (unlikely to overfit, 訓練集數量遠大於特徵數量)

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

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

\* 一定需有足夠特徵以提供充足的資訊來預測  $y$ , 否則多好的算法對預測品質也不會有太大幫助

人類專家看到這些特徵，有辦法預測出  $y$  嗎？

以及是否有足夠的訓練集