

## Building a spam classifier

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

- Collect lots of data
  - E.g. “honeypot” project.
- Develop sophisticated features based on email routing information (from email header).
- Develop sophisticated features for message body, e.g. should “discount” and “discounts” be treated as the same word? How about “deal” and “Dealer”? Features about punctuation?
- Develop sophisticated algorithm to detect misspellings (e.g. m0rtgage, med1cine, w4tches.)

误差分析

## 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 examples it is making errors on.

## Error Analysis

$m_{CV} = 500$  examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

- (i) What type of email it is *pharma, replica, steal passwords, ...*
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma: *12*

Replica/fake: *4*

➤ Steal passwords: *53*

Other: *31*

→ Deliberate misspellings: *5*  
(m0rgage, med1cine, etc.)

→ Unusual email routing: *16*

→ Unusual (spamming) punctuation: *32*

## The importance of numerical evaluation

Should discount/discounts/discounted/discounting be treated as the same word?

Can use “stemming” software (E.g. “Porter stemmer”)

universe/university.

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.

Need numerical evaluation (e.g., cross validation error) of algorithm’s performance with and without stemming.

Without stemming: *5% error* With stemming: *3% error*

Distinguish upper vs. lower case (Mom/mom): *3.2%*

偏斜类

## Cancer classification example

Train logistic regression model  $h_{\theta}(x)$ . ( $y = 1$  if cancer,  $y = 0$  otherwise)

Find that you got 1% error on test set.  
(99% correct diagnoses)

Only 0.50% of patients have cancer.

skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

0.5% error

99.2% accuracy (0.8% error)

降低至0.5% → 99.5% accuracy (0.5% error)

查准率 precision 召回率 recall

## Precision/Recall

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

Actual class

Predicted class	1	0
1	True positive	False positive
0	False negative	True negative

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

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

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

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

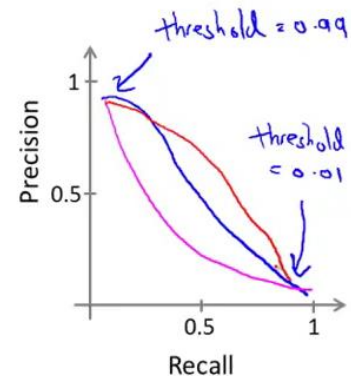
不同 threshold 对应不同 precision/recall 值，用 F1 值进行评估选择

## Trading off precision and recall

- Logistic regression:  $0 \leq h_{\theta}(x) \leq 1$   
 Predict 1 if  $h_{\theta}(x) \geq 0.5$  ~~0.5~~ ~~0.7~~ ~~0.9~~ ~~0.3~~ ←  
 Predict 0 if  $h_{\theta}(x) < 0.5$  ~~0.5~~ ~~0.7~~ ~~0.9~~ ~~0.3~~
- Suppose we want to predict  $y = 1$  (cancer) only if very confident.  
 → Higher precision, lower recall.
- Suppose we want to avoid missing too many cases of cancer (avoid false negatives).  
 → Higher recall, lower precision.

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$

$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



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

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

How to compare precision/recall numbers?

	Precision(P)	Recall (R)	<del>Average</del>	F <sub>1</sub> Score
→ Algorithm 1	<u>0.5</u>	<u>0.4</u>	<del>0.45</del>	0.444 ←
→ Algorithm 2	<u>0.7</u>	<u>0.1</u>	<del>0.4</del>	0.175 ←
Algorithm 3	<u>0.02</u>	1.0	<del>0.51</del>	0.0392 ←

Average:  ~~$\frac{P+R}{2}$~~

F<sub>1</sub> Score:  $2 \frac{PR}{P+R}$

$P=0$  or  $R=0 \Rightarrow F\text{-score} = 0.$   
 $P=1$  and  $R=1 \Rightarrow F\text{-score} = 1$

Predict  $y=1$  all the time