

Building a spam classifier

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

Features x : Choose 100 words indicative of spam/not spam.

E.g. deal, buy, discount, andrew, now, ...

$$X = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \begin{matrix} \text{andrew} \\ \text{buy} \\ \text{deal} \\ \text{discount} \\ \vdots \\ \text{now} \\ \vdots \end{matrix}$$

$$x \in \mathbb{R}^{100}$$

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise} \end{cases}$$

From: cheapsales@buystufffromme.com
To: ang@cs.stanford.edu
Subject: Buy now!

Deal of the week! Buy now!

Note: In practice, take most frequently occurring n words (10,000 to 50,000) in training set, rather than manually pick 100 words.

然后把这些单词

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.)

发垃圾邮件的也很机智 他们这么做就逃避了一些过滤

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*

更加复杂的特征变量 (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%*

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.

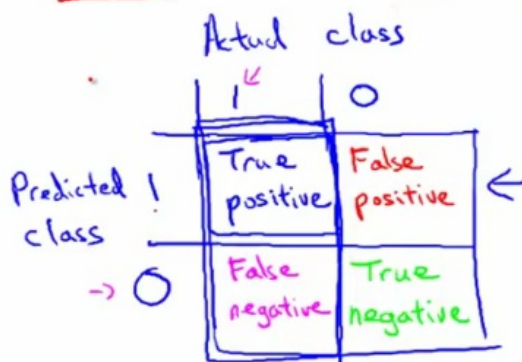
Without stemming: 5% error With stemming: 3% error

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

查准率 (Precision) 和召回率 (Recall) :

Precision/Recall

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



→ Precision

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

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

→ Recall

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

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

出现较少的类

查准率和召回率的折中 :

Trading off precision and recall

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

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

→ 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~~ ←

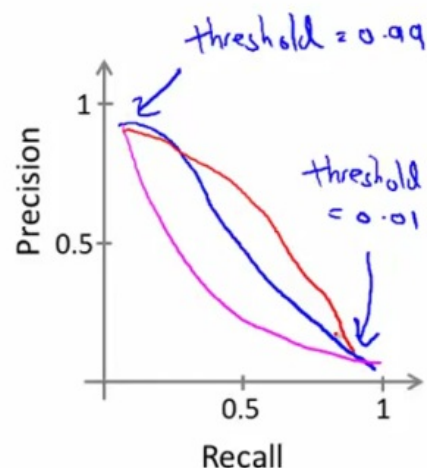
Predict 0 if $h_{\theta}(x) < 0.5$ ~~0.5~~ ~~0.7~~ ~~0.9~~ ~~0.3~~

→ 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.



这取决于回归模型的具体算法

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

自动选择临界值？

F₁ Score (F score)

How to compare precision/recall numbers?

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

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

Predict $y=1$ all the time

F₁ Score: $2 \frac{PR}{P+R}$

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

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

机器学习的数据问题：

数据的重要性：

Designing a high accuracy learning system

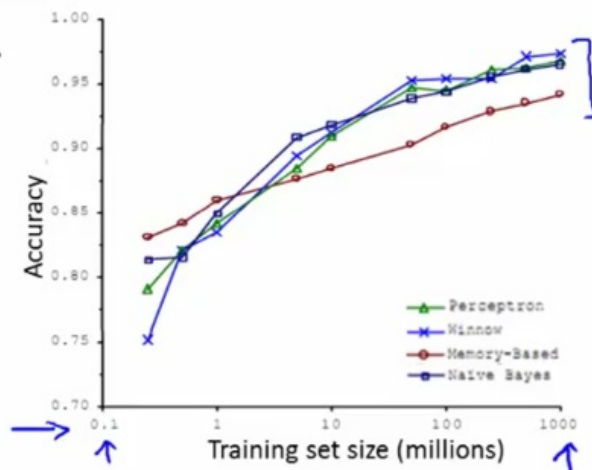
E.g. Classify between confusable words.

{to, two, too} {then, than}

→ For breakfast I ate two eggs.

Algorithms

- - Perceptron (Logistic regression)
- - Winnow
- - Memory-based
- - Naïve Bayes



"It's not who has the best algorithm that wins.

It's who has the most data."

的最佳方式 而不是

Banko and Brill, 2001]

Large data rationale

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

Example: For breakfast I ate two eggs.

Counterexample: Predict housing price from only size (feet²) and no other features.

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

Large data rationale

- 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

很好地测试测试数据集