# Week 6 — part 2

> **Machine Learning System design: Prioritizing what to work on: Spam classification example**

Supervised learning. $x$ = features of email

$$y = \text{spam } (1)$$
$$\text{not spam } (0)$$

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

e.g  deal, buy, discount__ = spam

Andrew, now... = not spam

$$x = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \\ \vdots \\ \vdots \end{bmatrix} \begin{matrix} \text{Andrew} \\ \text{Buy} \\ \text{deal} \\ \text{discount} \\ \vdots \\ \text{now} \\ \vdots \end{matrix}$$

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

example email

from:
To:
subject: buy now!

<u>Deal</u> of the week!

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

<u>Note</u>: In practise, take most frequency occurring $n$ words (10 000 to 50000) in training set, rather than manually pick 100 words.

# 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 "discounts" and "discount" be treated as the same word.
- develop sophisticated algorithms to detect misspellings

# Week 6

## Machine learning: Error Analysis

### Recommended approach

- Start with a simple algorithm that you can implement quickly. Implement it and test it on your cross-validation error data.
  └─ "Quick and dirty"
- 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 on in what type of examples it is making errors on.

└─ ways of avoiding "premature-optimization" //
└─ let evidence decide on where to spend our time rather than use gut feeling.

## 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, replicas, steal passwords.
   (ii) What cues (features) you think would have helped the algorithm classify correctly.

Pharma: 12
Replical fake: 4
steal passwords 53
other: 31

deliberate misspellings 5
(m6tgage, med1cine etc)
Unusual email routing: 16
─○ Unusual (spamming) punctuation: (32)
                                      ↗ strong signal, maybe worth the time working on this problem.

# Week 6

## Machine learning: error analysis

### 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) - potential error here.

error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try 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% error

if we use this as another parameter, by evaluating the error that the algorithm produces, we can see if the choice of the feature was correct.

Recommend implementing a "quick & dirty" implementation of your learning algorithm.

## Machine learning system design - error metrics for skewed classes

### 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

→ I.e hardly ever 1 (the patient has cancer) ∴ data is "skewed"

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

→ 0.5% error?

99.2% accuracy (0.8% error)
99.5% accuracy (0.5% error)

### Precision / Recall

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

|  | Actual class | |
|---|---|---|
|  | 1 | 0 |
| Predicted Class  1 | True positive | False positive |
| Predicted Class  0 | False Negative | True negative |

$y=0$
recall $=0$ ← not a good classifier

Computing precision and recall will give us a good sense of how well our classifier is doing.

**Precision**
(of all the patients that we predicted $y=1$, what fraction actually has cancer?)

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

divid by

**Recall**
(of all patients that actually have cancer, what fraction did we correctly detect having cancer?)

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

high recall is a good thing

# Week 6

## Machine learning: trading off precision and recall

Trading off precision and recall

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

Predict 1 if $\quad h_\theta(x) \geq 0.5 \,|\, 0.7 \,|\, 0.9 \,^{0.3}$

Predict 0 if $\quad h_\theta(x) < 0.5 \,|\, 0.7 \,|\, 0.9 \,|\, 0.3$

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

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

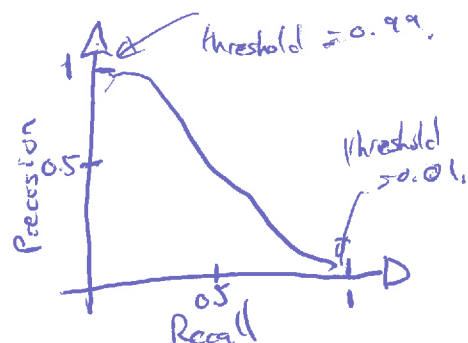Suppose we want to predict $y=1$ (cancer)
only if very confident.
→ higher precision, (lower recall)?  → one way of doing this

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

higher recall, lower precision

more generally: Predict 1 if $h_\theta(x) \geq$ threshold.


threshold = 0.99.
threshold = 0.01.
Precision (y-axis), Recall (x-axis)

## $F_1$ 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: $\dfrac{P+R}{2}$ ← this is one way

$F_1$ score: $2 \dfrac{PR}{P+R}$ → $P=0$ or $R=0$
$\Rightarrow$ Fscore $= 0$
$P=1$ & $R=1 \Rightarrow$ Fscore $=1$

Predict $y=1$ all the time → not a good algorithm get high average.

→ averaging is not a good method of evaluating the algorithm's overall performance.

# Week 6

Machine learning system design: data for machine learning

## Design a high accuracy learning system

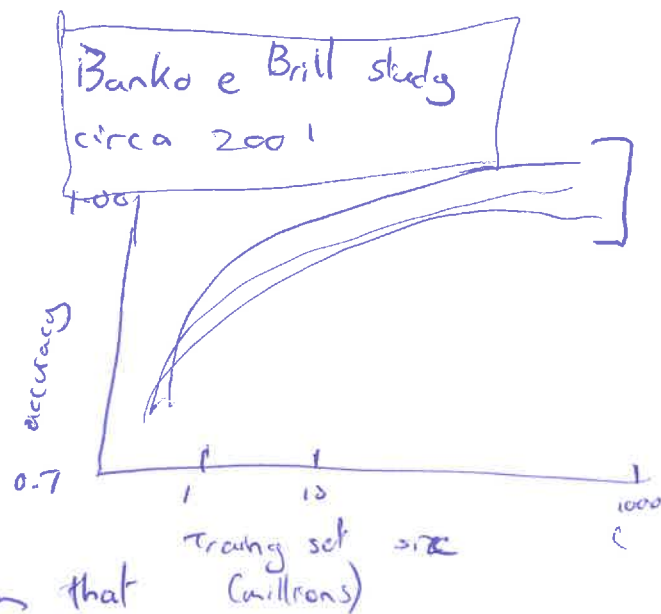e.g. Classify between confusable words.
{to, two, too} / {then, than}
For breakfast I ate __two__ eggs.

Algorithms
- Perceptron (logistic regression)
- Winnow
- Memory based
- Naïve based.



Banko e Brill study circa 2001

"It's not who has the best algorithm that wins,
it's who has the most data"

## 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

Counter example: Predict housing price from only size (feet$^2$) and no other features. → This is an example where it would be very difficult to predict the price accurately.

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

↳ So ask a human expert if the chosen parameter is 'good enough' to make the prediction.

# 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.

→ $J_{train}(\theta)$ will be small.

Use a very large training set (unlikely to overfit) ← because we have such a massive trng set.

low variance

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

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