

## CS1110 lec 12: Analysis of lies using recursion 07 Oct 2010

Today: recursion in an interdisciplinary application:  
computer science/computational linguistics, psychology, history/politics

Have your iClickers out.

Reading for next lecture (casting about): Secs 4.2 & 4.3

- *Prelim, 7:30-9pm today*
  - Last name A-K: go to Olin 155
  - Last name L-Z: go to Olin 255
- A4 due Saturday Oct. 16
- No labs next week (Tue-Wed Oct 11-12), due to fall break
- No office/consulting hours Friday through Tuesday inclusive (Oct 8-12), due to fall break

1

## Lies, damned lies, and statistics

James Pennebaker et al., "Lying words: predicting deception from linguistic styles", 2003:

**Claim:** deceptive communication is characterized by (among others):

- fewer 1<sup>st</sup>-person singular and 3<sup>rd</sup>-person pronouns ("I", "they")
- more negative emotion words ("hate", "enemy")
- fewer "complex/exclusive" words ("but", "except", "without")

**Research question (1):** What really are the best cues? (or models)

More realistic, convenient source of "lies"?

The "Iraq War Card False Statements Database"

<http://projects.publicintegrity.org/WarCard/Search/Default.aspx>

2

## Where we are, and why [besides automatic lie detection being inherently cool]

**Research question 1:** What are the best linguistic lie cues? (or models)

**Research question 2:** given statements regarding Iraq by top Bush administration officials, are the "true" and "false" ones distinguishable? This would imply something about their beliefs.

- Demonstration of interdisciplinary research involving computer science, psychology/linguistics, politics and history
- Demonstration of methodology in approaching a programming problem
  - stepwise refinement, writing and reading specs carefully, String manipulation, recursion, testing, etc.

(Lecture loosely based on joint work with CS grad student Cristian Danescu-Niculescu-Mizil and CS undergrad Haden Lee, in consultation with Comm. Prof. Jeff Hancock.)

3

## Formulating the task: might something be a lie cue?

### Given

- a target word  $w$  (e.g., "they")
- a file containing "lie span mark-up":

... with respect to Iraq, the problem is quite simple. **<font color="red">**we suspect **they** are developing weapons of mass destruction. we **more** than suspect it; we know it. **</font>** **they** could...

spanStart                      spanEnd

**Consider** these two statistics for  $w$ :

- number of **hits**: occurrences of  $w$  in a "lie" span
- number of **misses**: occurrences of  $w$  not in a "lie" span

**Write** a class `LieData` with method `counts(w,...)` that will tell us the number of hits and the number of misses for  $w$  in a text (file).

4

To write: a class `LieData` with method `counts(w)`.  
We need to track the source file's text, `spanStart`, `spanEnd`, and the target word  $w$ .

### Q: How should we declare the relevant entities?

(A) *text and the span delimiters are stored in individual objects*  
`private String text; private String spanStart; private String spanEnd;`  
`public String counts(String w) {...}`

(B) *everything is a parameter to static method counts*  
`public static String counts(String textFileName, String spanStart,`  
`String spanEnd, String w) {...}`

(C) *this mixture of objects and static:*  
`private String text; private String spanStart; private String spanEnd`  
`public static counts(String w) {...}`

5

1. Get data from source file into Java-manageable format  
create new object via "new `LieData(<file>)`" with String field `text`
2. Get target word  $w$   
parameter for method `counts(w)`
3. Process each occurrence of  $w$  in the source file's text
  - Is it in or not in a span?
4. Report relevant statistics  
output of `counts(w)`: "hits <h>; misses <m>"

6

**The core recursive idea: how many hits past an index  $i$ ?**

$w$ : the target word, i.e., "they"  
 $i$ : a starting index

$iW$ : index of 1<sup>st</sup> occurrence of  $w$  after index  $i$   
 $iW + w.length()$

`text: "... end of it all they will bring to the table a new means..."`

if  $h$  is the number of hits for  $w$  here (i.e., past  $iW + w.length()$ ) ...

... then if the "they" at  $iW$  is a hit, the # of hits past  $i$  is  $h + 1$ ;  
 if the "they" at  $iW$  is a miss, the # of hits past  $i$  is  $h$ .

Handy version of `indexOf`: `iW = text.indexOf(w, i)`.

7

We need the output of helper method `countsFrom(String w, int iW)` to encode *both* the numbers of hits *and* the number of misses in a way such that we can extract or add to them.

**Q: What format should `countsFrom`'s output be?**

- (A) an int, the difference between the # of hits and the # of misses
- (B) a String "hits: <# of hits>; misses: <#misses>"
- (C) a double, using decimal as separator. Example: 4.5 means 4 hits, 5 misses
- (D) an object that maintains hits and misses (we would have to write the class)

8

**A helper class for storing hit/miss pairs**

```

/** An instance maintains a number of hits and a number
of misses */
public class TallyPair {
    public int hits; // number of hits
    public int misses; // number of misses

    /** Constructor: a TallyPair with 0 hits and 0 misses */
    public TallyPair() {
        ;
    }
    /** = hits: <# of hits>; misses <# of misses> */
    public String toString() {
        return "hits: " + hits + "; misses: " + misses;
    }
}

```

(A better choice would be to make this a private static nested class: see pp 348–350.)

9

**How do we determine if  $w$  at index  $iW$  is a hit or not?**

**boolean procedure `isHit(String w, int iW)`**

... they ... (no spanEnds) ... spanEnd ...

$iW$ , an index of  $w$        $iEnd$ : index of closest *spanEnd* after  $iW$

Handy version of `lastIndexOf`:  
**return** `iEnd != -1 && text.lastIndexOf(spanStart, iEnd) < iW`;

10