**CS276 – Information Retrieval and Web Search**

Checking in. By the end of this week you need to have:
- Watched the online videos corresponding to the first 6 chapters of *IIR* **or/and** read chapters 1–6 of the book
- Done programming assignment 1 (due Thursday)
- Submitted 5 search queries for the Stanford domain (for PA3)
- Oh, and problem set 1 was due last Thursday ☺

Today: Probabilistic models of spelling correction for PA2
- You should also look at chapter 3 video/book for other material

Thursday: Class lab on map-reduce

---

# Spelling Correction and the Noisy Channel

## The Spelling Correction Task

---

## Applications for spelling correction



Word processing

Phones

Web search

3

---

## Spelling Tasks

- Spelling Error Detection
- Spelling Error Correction:
  - Autocorrect
    - hte→the
  - Suggest a correction
  - Suggestion lists

4

---

## Types of spelling errors

- Non-word Errors
  - *graffe* →*giraffe*
- Real-word Errors
  - Typographical errors
    - *three* →*there*
  - Cognitive Errors (homophones)
    - *piece*→*peace*,
    - *too* → *two*

5

---

## Rates of spelling errors

**26**%:  Web queries  Wang *et al.* 2003

**13**%:  Retyping, no backspace: Whitelaw *et al.* English&German

**7**%: Words corrected retyping on phone-sized organizer

**2**%: Words uncorrected on organizer Soukoreff &MacKenzie 2003

**1-2**%:  Retyping: Kane and Wobbrock 2007, Gruden et al. 1983

6

---

## Non-word spelling errors

- Non-word spelling error detection:
  - Any word not in a **dictionary** is an error
  - The larger the dictionary the better
- Non-word spelling error correction:
  - Generate **candidates**: real words that are similar to error
  - Choose the one which is best:
    - Shortest weighted edit distance
    - Highest noisy channel probability

7

## Real word spelling errors

- For each word *w*, generate candidate set:
  - Find candidate words with similar **pronunciations**
  - Find candidate words with similar **spelling**
  - Include *w* in candidate set
- Choose best candidate
  - Noisy Channel

8

# Spelling Correction and the Noisy Channel

### The Noisy Channel Model of Spelling

## Noisy Channel Intuition



10

## Noisy Channel aka Bayes' Rule

- We see an observation *x* of a misspelled word
- Find the correct word $\hat{w}$

$$\hat{w} = \underset{w \in V}{\operatorname{argmax}} P(w \mid x)$$

$$= \underset{w \in V}{\operatorname{argmax}} \frac{P(x \mid w)P(w)}{P(x)}$$

$$= \underset{w \in V}{\operatorname{argmax}} P(x \mid w)P(w)$$

11

## History: Noisy channel for spelling proposed around 1990

- **IBM**
  - Mays, Eric, Fred J. Damerau and Robert L. Mercer. 1991. Context based spelling correction. *Information Processing and Management*, 23(5), 517–522
- **AT&T Bell Labs**
  - Kernighan, Mark D., Kenneth W. Church, and William A. Gale. 1990. A spelling correction program based on a noisy channel model. Proceedings of COLING 1990, 205-210

2

### Non-word spelling error example

acress

13

### Candidate generation

- Words with similar spelling
  - Small edit distance to error
- Words with similar pronunciation
  - Small edit distance of pronunciation to error

14

### Damerau-Levenshtein edit distance

- Minimal edit distance between two strings, where edits are:
  - Insertion
  - Deletion
  - Substitution
  - Transposition of two adjacent letters

- See *IIR* sec 3.3.3 for edit distance

15

### Words within 1 of acress

| Error | Candidate Correction | Correct Letter | Error Letter | Type |
|-------|---------------------|----------------|--------------|------|
| acress | actress | t | – | deletion |
| acress | cress | – | a | insertion |
| acress | caress | ca | ac | transposition |
| acress | access | c | r | substitution |
| acress | across | o | e | substitution |
| acress | acres | – | s | insertion |
| acress | acres | – | s | insertion |

16

### Candidate generation

- 80% of errors are within edit distance 1
- Almost all errors within edit distance 2

- Also allow insertion of **space** or **hyphen**
  - thisidea → this idea
  - inlaw → in-law

17

### Wait, how do you generate the candidates?

1. Run through dictionary, check edit distance with each word
2. Generate all words within edit distance ≤ $k$ (e.g., $k$ = 1 or 2) and then intersect them with dictionary
3. Use a character $k$-gram index and find dictionary words that share "most" $k$-grams with word (e.g., by Jaccard coefficient)
   - see *IIR* sec 3.3.4
4. Compute them fast with a Levenshtein finite state transducer
5. Have a precomputed hash of words to possible corrections

18

3

## Language Model

- Just use the unigram probability of words
  - Take big supply of words (your document collection with $T$ tokens)

$$P(w) = \frac{C(w)}{T}$$

19

## Unigram Prior probability

Counts from 404,253,213 words in Corpus of Contemporary English (COCA)

| word | Frequency of word | P(word) |
|------|------|------|
| actress | 9,321 | .0000230573 |
| cress | 220 | .0000005442 |
| caress | 686 | .0000016969 |
| access | 37,038 | .0000916207 |
| across | 120,844 | .0002989314 |
| acres | 12,874 | .0000318463 |

20

## Channel model probability

- **Error model probability, Edit probability**
- *Kernighan, Church, Gale 1990*

- *Misspelled word $x = x_1, x_2, x_3... x_m$*
- *Correct word $w = w_1, w_2, w_3..., w_n$*

- P(x|w) = probability of the edit
  - (deletion/insertion/substitution/transposition)

21

## Computing error probability: confusion matrix

```
del[x,y]:     count(xy typed as x)
ins[x,y]:     count(x typed as xy)
sub[x,y]:     count(x typed as y)
trans[x,y]:   count(xy typed as yx)
```

Insertion and deletion conditioned on previous character

22

## Confusion matrix for spelling errors



**sub[X, Y] = Substitution of X (incorrect) for Y (correct)**

## Generating the confusion matrix

- Peter Norvig's list of errors
- Peter Norvig's list of counts of single-edit errors

  - All Peter Norvig's ngrams data links: http://norvig.com/ngrams/

24

### Slide 25

Christopher Manning

## Channel model

Kernighan, Church, Gale 1990

$$P(x|w) = \begin{cases} \dfrac{\text{del}[w_{i-1},w_i]}{\text{count}[w_{i-1}w_i]}, & \text{if deletion} \\[2mm] \dfrac{\text{ins}[w_{i-1},x_i]}{\text{count}[w_{i-1}]}, & \text{if insertion} \\[2mm] \dfrac{\text{sub}[x_i,w_i]}{\text{count}[w_i]}, & \text{if substitution} \\[2mm] \dfrac{\text{trans}[w_i,w_{i+1}]}{\text{count}[w_iw_{i+1}]}, & \text{if transposition} \end{cases}$$

25

### Slide 26

Christopher Manning

## Smoothing probabilities: Add-1 smoothing

- But if we use the last slide, unseen errors are impossible!
- They'll make the overall probability 0. That seems too harsh
  - e.g., in Kernighan's chart q➜a and a➜q are both 0, even though they're adjacent on the keyboard!
- A simple solution is to add one to all counts and then if there is a |A| character alphabet, to normalize appropriately:

$$\text{If substitution, } P(x\,|\,w) = \frac{\text{sub}[x,w]+1}{\text{count}[w]+A}$$

26

### Slide 27

Christopher Manning

## Channel model for `acress`

| Candidate Correction | Correct Letter | Error Letter | x\|w | P(x\|word) |
|---|---|---|---|---|
| actress | t | – | c\|ct | .000117 |
| cress | – | a | a\|# | .00000144 |
| caress | ca | ac | ac\|ca | .00000164 |
| access | c | r | r\|c | .000000209 |
| across | o | e | e\|o | .0000093 |
| acres | – | s | es\|e | .0000321 |
| acres | – | s | ss\|s | .0000342 |

27

### Slide 28

Christopher Manning

## Noisy channel probability for `acress`

| Candidate Correction | Correct Letter | Error Letter | x\|w | P(x\|word) | P(word) | $10^9$ *P(x\|w)P(w) |
|---|---|---|---|---|---|---|
| actress | t | – | c\|ct | .000117 | .0000231 | 2.7 |
| cress | – | a | a\|# | .00000144 | .000000544 | .00078 |
| caress | ca | ac | ac\|ca | .00000164 | .00000170 | .0028 |
| access | c | r | r\|c | .000000209 | .0000916 | .019 |
| across | o | e | e\|o | .0000093 | .000299 | 2.8 |
| acres | – | s | es\|e | .0000321 | .0000318 | 1.0 |
| acres | – | s | ss\|s | .0000342 | .0000318 | 1.0 |

28

### Slide 29

Christopher Manning

## Noisy channel probability for `acress`

| Candidate Correction | Correct Letter | Error Letter | x\|w | P(x\|word) | P(word) | $10^9$ *P(x\|w)P(w) |
|---|---|---|---|---|---|---|
| actress | t | – | c\|ct | .000117 | .0000231 | 2.7 |
| cress | – | a | a\|# | .00000144 | .000000544 | .00078 |
| caress | ca | ac | ac\|ca | .00000164 | .00000170 | .0028 |
| access | c | r | r\|c | .000000209 | .0000916 | .019 |
| **across** | **o** | **e** | **e\|o** | **.0000093** | **.000299** | **2.8** |
| acres | – | s | es\|e | .0000321 | .0000318 | 1.0 |
| acres | – | s | ss\|s | .0000342 | .0000318 | 1.0 |

29

### Slide 30

Christopher Manning

## Incorporating context words: Context-sensitive spelling correction

- Determining whether **actress** or **across** is appropriate will require looking at the context of use
- We can do this with a better **language model**
  - You learned/can learn a lot about language models in CS124 or CS224N
  - Here we present just enough to be dangerous/do the assignment
- A **bigram language model** conditions the probability of a word on (just) the previous word

  $P(w_1...w_n) = P(w_1)P(w_2|w_1)...P(w_n|w_{n-1})$

30

## Incorporating context words

- For unigram counts, P($w$) is always non-zero
  - if our dictionary is derived from the document collection
- This won't be true of P($w_k|w_{k-1}$). We need to **smooth**
- We could use add-1 smoothing on this conditional distribution
- But here's a better way: interpolate a unigram and a bigram:

  $P_{li}(w_k|w_{k-1}) = \lambda P_{uni}(w_1) + (1-\lambda)P_{mle}(w_k|w_{k-1})$
  - $P_{mle}(w_k|w_{k-1}) = C(w_k|w_{k-1}) / C(w_{k-1})$
  - This is called a "maximum likelihood estimate" (mle)
  - For categorical variables you get an mle by just counting and dividing

31

---

## All the important fine points

- Our unigram probability $P_{uni}(w_k) = C(w_k) / T$ is also an mle
  - This is okay if our dictionary is only words in the document collection – will be non-zero
  - Otherwise we'd need to smooth it to avoid zeroes (e.g., add-1 smoothing)
- Note that we have several probability distributions for words
  - Keep them straight!
- You might want/need to work with log probabilities:
  - $\log P(w_1 \ldots w_n) = \log P(w_1) + \log P(w_2|w_1) + \ldots + \log P(w_n|w_{n-1})$
  - Otherwise, be very careful about floating point underflow
- Our query may be words anywhere in a document
  - We'll start the bigram estimate of a sequence with a unigram estimate
  - Often, people instead condition on a start-of-sequence symbol, but not good here
  - Because of this, the unigram and bigram counts have different totals. Not a problem

32

---

## Using a bigram language model

- "a stellar and versatile **acress** whose combination of sass and glamour…"
- Counts from the Corpus of Contemporary American English with add-1 smoothing
- P(actress|versatile)=.000021 P(whose|actress) = .0010
- P(across|versatile) =.000021 P(whose|across) = .000006

- P("versatile actress whose") = .000021*.0010 = 210 x10$^{-10}$
- P("versatile across whose")  = .000021*.000006 = 1 x10$^{-10}$

33

---

## Using a bigram language model

- "a stellar and versatile **acress** whose combination of sass and glamour…"
- Counts from the Corpus of Contemporary American English with add-1 smoothing
- P(actress|versatile)=.000021 P(whose|actress) = .0010
- P(across|versatile) =.000021 P(whose|across) = .000006

- **P("versatile actress whose") = .000021*.0010 = 210 x10$^{-10}$**
- P("versatile across whose")  = .000021*.000006 = 1 x10$^{-10}$

34

---

## Evaluation

- Some spelling error test sets
  - Wikipedia's list of common English misspelling
  - Aspell filtered version of that list
  - Birkbeck spelling error corpus
  - Peter Norvig's list of errors (includes Wikipedia and Birkbeck, for training or testing)

35

---

# Spelling Correction and the Noisy Channel

## Real-Word Spelling Correction

## Real-word spelling errors

- …leaving in about fifteen **_minuets_** to go to her house.
- The design **_an_** construction of the system…
- Can they **_lave_** him my messages?
- The study was conducted mainly **_be_** John Black.

- 25-40% of spelling errors are real words    Kukich 1992

37

---

## Solving real-word spelling errors

- For each word in sentence
  - Generate *candidate set*
    - the word itself
    - all single-letter edits that are English words
    - words that are homophones
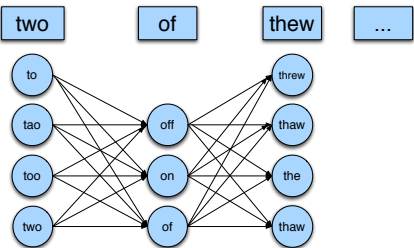- Choose best candidates
  - Noisy channel model

38

---

## Noisy channel for real-word spell correction

- Given a sentence $w_1, w_2, w_3, …, w_n$
- Generate a set of candidates for each word $w_i$
  - Candidate($w_1$) = {$w_1, w'_1, w''_1, w'''_1, …$}
  - Candidate($w_2$) = {$w_2, w'_2, w''_2, w'''_2, …$}
  - Candidate($w_n$) = {$w_n, w'_n, w''_n, w'''_n, …$}
- Choose the sequence W that maximizes P(W)
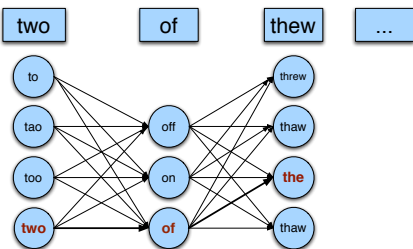
---

## Noisy channel for real-word spell correction



40

---

## Noisy channel for real-word spell correction



41

---

## Simplification: One error per sentence

- Out of all possible sentences with one word replaced
  - $w_1$, **$w''_2$**, $w_3$, $w_4$        two **off** thew
  - $w_1$, $w_2$, **$w'_3$**, $w_4$        two of **the**
  - **$w'''_1$**, $w_2$, $w_3$, $w_4$        **too** of thew
  - …
- Choose the sequence W that maximizes P(W)

---

7

## Where to get the probabilities

- Language model
  - Unigram
  - Bigram
  - etc.
- Channel model
  - Same as for non-word spelling correction
  - Plus need probability for no error, P(w|w)

43

## Probability of no error

- What is the channel probability for a correctly typed word?
- P("the"|"the")
  - If you have a big corpus, you can estimate this percent correct

- But this value depends strongly on the application
  - .90 (1 error in 10 words)
  - .95 (1 error in 20 words)
  - .99 (1 error in 100 words)

44

## Peter Norvig's "thew" example

| x | w | x\|w | P(x\|w) | P(w) | $10^9$ P(x\|w)P(w) |
|---|---|------|---------|------|-------------------|
| thew | the | ew\|e | 0.000007 | 0.02 | 144 |
| thew | thew | | 0.95 | 0.00000009 | 90 |
| thew | thaw | e\|a | 0.001 | 0.0000007 | 0.7 |
| thew | threw | h\|hr | 0.000008 | 0.000004 | 0.03 |
| thew | thwe | ew\|we | 0.000003 | 0.00000004 | 0.0001 |

45

## State of the art noisy channel

- We never just multiply the prior and the error model
- Independence assumptions→probabilities not commensurate
- Instead: Weight them

$$\hat{w} = \underset{w \in V}{\mathrm{argmax}}\, P(x \mid w) P(w)^{\lambda}$$

- Learn λ from a development test set

46

## Improvements to channel model

- Allow richer edits   (Brill and Moore 2000)
  - ent→ant
  - ph→f
  - le→al
- Incorporate pronunciation into channel (Toutanova and Moore 2002)
- Incorporate device into channel

47

## Nearby keys