

UNIVERSITY of DELAWARE

Spelling Correction

(Search and Data Mining)

Hui Fang
Department of Electrical and Computer Engineering
University of Delaware

1

UNIVERSITY of DELAWARE

Reference

- How to write a spelling corrector, by Peter Norvig
<http://norvig.com/spell-correct.html>

UNIVERSITY of DELAWARE

Types of Spelling Errors

- Errors are not real words.
– *elephant* → *elephant*
- Errors are real words.
• *three* → *there*
• *their* → *there*

3

UNIVERSITY of DELAWARE

How to Detect and Correct Non-word Spelling Errors

- Check whether the word is in a dictionary.
- If no, need to generate correction
 - Generate candidates
 - Choose the one which is best

4

UNIVERSITY of DELAWARE

Candidate Generation

- Find real-words that are similar to the errors
 - Similar spelling
 - Similar pronunciation
- Use minimal edit distance between two strings to measure the similarity
 - Edits include insertion, deletion, substitution and transposition of two adjacent letters.
 - 80% of errors are within edit distance 1
 - Almost all errors within edit distance 2

5

UNIVERSITY of DELAWARE

Ranking Candidates

- Given a misspelled word x , we can rank possible correct words w based on

$$\hat{w} = \operatorname{argmax}_{w \in \text{candidates}} P(w | x)$$

$$= \operatorname{argmax}_{w \in \text{candidates}} \frac{P(x | w)P(w)}{P(x)}$$

$$= \operatorname{argmax}_{w \in \text{candidates}} P(x | w)P(w)$$

7

UNIVERSITY of DELAWARE

Ranking Candidates

- Given a misspelled word x , we can rank possible correct words w based on

$$\begin{aligned}\hat{w} &= \operatorname{argmax}_{w \in \text{candidates}} P(w|x) \\ &= \operatorname{argmax}_{w \in \text{candidates}} \frac{P(x|w)P(w)}{P(x)} \\ &= \operatorname{argmax}_{w \in \text{candidates}} P(x|w)P(w)\end{aligned}$$

Go through all candidates, and choose the one that gives the best combined probability score

8

UNIVERSITY of DELAWARE

Ranking Candidates

- Given a misspelled word x , we can rank possible correct words w based on

$$\begin{aligned}\hat{w} &= \operatorname{argmax}_{w \in \text{candidates}} P(w|x) \\ &= \operatorname{argmax}_{w \in \text{candidates}} P(x|w)P(w)\end{aligned}$$

How likely do you see w in a general collection?

Can be estimated using any language models, such as unigram models.

9

UNIVERSITY of DELAWARE

Ranking Candidates

- Given a misspelled word x , we can rank possible correct words w based on

$$\begin{aligned}\hat{w} &= \operatorname{argmax}_{w \in V} P(w|x) \\ &= \operatorname{argmax}_{w \in V} \frac{P(x|w)P(w)}{P(x)} \\ &= \operatorname{argmax}_{w \in V} P(x|w)P(w)\end{aligned}$$

You want to type x , how likely is it to type x instead?

10

UNIVERSITY of DELAWARE

Error Model : $P(x|w)$

- $P(x|w)$ = probability of the edit
 - deletion[x,y]: count(xy typed as x)
 - insertion[x,y]: count(x typed as xy)
 - substitution[x,y]: count(x typed as y)
 - transposition[x,y]: count(xy typed as yx)

It requires lots of training data!

11

UNIVERSITY of DELAWARE

A simplified error model : $P(x|w)$

- A simplified solution without training data set
 - All known words of edit distance 1 are infinitely more probable than known words of edit distance 2.
 - All known words of edit distance 1 are infinitely less probable than a known word of edit distance 0.

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

12

UNIVERSITY of DELAWARE

Peter Norvig's Spelling Corrector

- Not industrial-strength spell corrector
- A toy non-word spelling corrector
 - 80 or 90% accuracy
 - At a processing speed of at least 10 words per second

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

14

```

import re, collections

def words(text): return re.findall('[a-z]+', text.lower())

def train(features):
    model = collections.defaultdict(lambda: 1)
    for f in features:
        model[f] += 1
    return model

NWORDS = train(words(file('big.txt').read()))

alphabet = 'abcdefghijklmnopqrstuvwxyz'

def edits1(word):
    splits   = [(word[:i], word[i:]) for i in range(len(word) + 1)]
    deletes  = [a + b[1:] for a, b in splits if b]
    transposes = [a + b[1] + b[0] + b[2:] for a, b in splits if len(b)>1]
    replaces = [a + c + b[1:] for a, b in splits for c in alphabet if b]
    inserts  = [a + c + b for a, b in splits for c in alphabet]
    return set(deletes + transposes + replaces + inserts)

def known_edits2(word):
    return set(e2 for e1 in edits1(word) for e2 in edits1(e1) if e2 in NWORDS)

def known(words): return set(w for w in words if w in NWORDS)

def correct(word):
    candidates = known([word]) or known(edits1(word)) or known_edits2(word) or {word}
    return max(candidates, key=NWORDS.get)

```

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

Recall: Candidate Generation

- Find real-words that are similar to the errors
 - Similar spelling
 - Similar pronunciation
- Use Edit distance to measure the similarity
 - 80% of errors are within edit distance 1
 - Almost all errors within edit distance 2

16

```

import re, collections

def words(text): return re.findall('[a-z]+', text.lower())

def train(features):
    model = collections.defaultdict(lambda: 1)
    for f in features:
        model[f] += 1
    return model

NWORDS = train(words(file('big.txt').read()))

alphabet = 'abcdefghijklmnopqrstuvwxyz'

def edits1(word):
    splits   = [(word[:i], word[i:]) for i in range(len(word) + 1)]
    deletes  = [a + b[1:] for a, b in splits if b]
    transposes = [a + b[1] + b[0] + b[2:] for a, b in splits if len(b)>1]
    replaces = [a + c + b[1:] for a, b in splits for c in alphabet if b]
    inserts  = [a + c + b for a, b in splits for c in alphabet]
    return set(deletes + transposes + replaces + inserts)

def known_edits2(word):
    return set(e2 for e1 in edits1(word) for e2 in edits1(e1) if e2 in NWORDS)

def known(words): return set(w for w in words if w in NWORDS)

def correct(word):
    candidates = known([word]) or known(edits1(word))
    return max(candidates, key=NWORDS.get)

    
```

Returns a set of all words w that are one edit away from x

Consider words that are two edits away.

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

Recall: Ranking Candidates

- Given a misspelled word x, we can rank possible correct words w based on

$$\hat{w} = \operatorname{argmax}_{w \in \text{candidates}} P(w|x)$$

How likely do you see w in a general collection?
- Can be estimated using any language models, such as unigram models.

$$\hat{w} = \operatorname{argmax}_{w \in \text{candidates}} P(x|w)P(w)$$

18

```

import re, collections

def words(text): return re.findall('[a-z]+', text.lower())

def train(features):
    model = collections.defaultdict(lambda: 1)
    for f in features:
        model[f] += 1
    return model

NWORDS = train(words(file('big.txt').read()))

alphabet = 'abcdefghijklmnopqrstuvwxyz'

def edits1(word):
    splits   = [(word[:i], word[i:]) for i in range(len(word) + 1)]
    deletes  = [a + b[1:] for a, b in splits if b]
    transposes = [a + b[1] + b[0] + b[2:] for a, b in splits if len(b)>1]
    replaces = [a + c + b[1:] for a, b in splits for c in alphabet if b]
    inserts  = [a + c + b for a, b in splits for c in alphabet]
    return set(deletes + transposes + replaces + inserts)

def known_edits2(word):
    return set(e2 for e1 in edits1(word) for e2 in edits1(e1))

def known(words): return set(w for w in words if w in NWORDS)

def correct(word):
    candidates = known([word]) or known(edits1(word)) or known_edits2(word) or {word}
    return max(candidates, key=NWORDS.get)

```

NWORDS[w] holds a count of how many times that word w has been seen in the training file. (add one smoothing)

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

Recall: Ranking Candidates

- Given a misspelled word x, we can rank possible correct words w based on

$$\hat{w} = \operatorname{argmax}_{w \in V} P(w|x)$$

You want to type x, how likely is it to type x instead?
- $$= \operatorname{argmax}_{w \in V} \frac{P(x|w)P(w)}{P(x)}$$
- $$= \operatorname{argmax}_{w \in V} P(x|w)P(w)$$

20

```

UNIVERSITY OF DELAWARE import re, collections
def words(text): return re.findall('[a-z]+', text.lower())
def train(features):
    model = collections.defaultdict(int)
    for f in features:
        model[f] += 1
    return model
NWORDS = train(words(file('big.txt')))
alphabet = 'abcdefghijklmnopqrstuvwxyz'
def edits1(word):
    splits   = [word[:i] + word[i+1:] for i in range(len(word))]
    deletes = [a + b[1:] for a in alphabet for b in splits]
    transposes = [a + b[1:] + b[0] + b[2:] for a in alphabet for b in splits]
    replaces = [a + c + b[1:] for a in alphabet for c in alphabet for b in splits]
    inserts = [a + c + b[1:] for a, b in zip(alphabet, splits) for c in alphabet]
    return set(deletes + transposes + replaces + inserts)
def known_edits2(word):
    return set(e2 for e1 in edits1(word) for e2 in edits1(e1) if e2 in NWORDS)
def known(words): return set(w for w in words if w in NWORDS)
def correct(word):
    candidates = known({word}) or known_edits1(word) or known_edits2(word) or {word}
    return max(candidates, key=NWORDS.get)

```

Without training set, defined a trivial model: all known words of edit distance 1 are infinitely more probable than known words of edit distance 2, and infinitely less probable than a known word of edit distance 0.

Reference: How to write a spelling corrector, by Peter Norvig, <http://norvig.com/spell-correct.html>

Solving real-word spelling errors

- Example:
 - It takes her about *6 minuets* to go back home.
- Methods:
 - For each word in sentence
 - Generate *candidate set*
 - the word itself
 - all single-letter edits that are English words
 - words that have the same pronunciation
 - Choose the most likely combination

23