
Noisy Channel

A decorative horizontal bar at the bottom of the slide, consisting of three segments: a red segment on the left, a light gray segment in the middle, and a medium gray segment on the right.

Noisy Channel

We see an observation x of the misspelled word

Find the correct word w

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

Noisy Channel

We see an observation x of the misspelled word

Find the correct word w

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

Noisy Channel

We see an observation x of the misspelled word

Find the correct word w

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

Creating candidate set

- Words with similar spelling
- Words with similar pronunciation

Some statistics

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

Allow deletion of space or hyphen

- thisidea → this idea
- inlaw → in-law

Words within edit distance 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

Computing error probability: confusion matrix

- $\text{del}[x,y]$: count (xy typed as x)
- $\text{ins}[x,y]$: count (x typed as xy)
- $\text{sub}[x,y]$: count (x typed as y)
- $\text{trans}[x,y]$: count(xy typed as yx)

Computing error probability: confusion matrix

- $\text{del}[x,y]$: count (xy typed as x)
- $\text{ins}[x,y]$: count (x typed as xy)
- $\text{sub}[x,y]$: count (x typed as y)
- $\text{trans}[x,y]$: count(xy typed as yx)

Insertion and deletion are conditioned on previous character

Channel model

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

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	.00000093
acres	-	s	es e	.0000321
acres	-	s	ss s	.0000342

Noisy channel probability for *acress*

Candidate Correction	Correct Letter	Error Letter	$x w$	$P(x \text{word})$	$P(\text{word})$	$10^9 * P(x w)P(w)$
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

- Here we have maximum probability 2.8 therefore correct candidate is **across**
- Without considering context, this method is very good.
- But suppose if we have probabilities values very near like 2.7, 2.71..., than without using context we cannot choose which candidate is correct candidate.
- **Example with context:**
- Consider the sentence “..... **Versatile across whose.....**”
- So here we have two candidates ‘**across**’ and ‘**actress**’
- $P(\text{actress} | \text{versatile}) = 0.000021$
- $P(\text{across} | \text{versatile}) = 0.000021$
- $P(\text{whose} | \text{actress}) = 0.0010$, $P(\text{whose} | \text{across}) = 0.000006$
- $P(\text{“versatile actress whose”}) = 0.000021 * 0.0010 = 210 * 10^{-10}$
- $P(\text{“versatile across whose”}) = 0.000021 * 0.000006 = 1 * 10^{-10}$
- Here **actress** is more probable.

Real-word spelling errors

- The study was conducted mainly **be** John Black
- The design **an** construction of the system ...

25-40% of spelling errors are real words

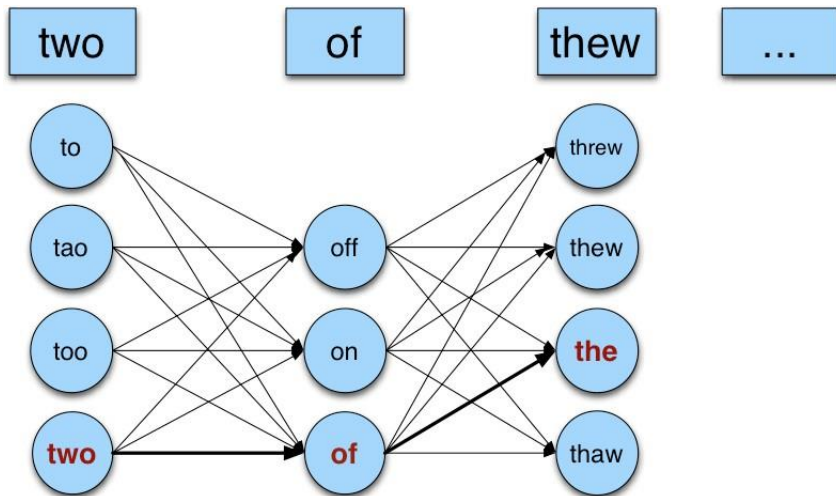
Noisy channel for real-word spell correction

Given a sentence $X = w_1, w_2, w_3 \dots, w_n$

- Candidate (w_1) = $\{w_1, w_1^j, w_1^{jj}, w_1^{jjj}, \dots\}$
- Candidate (w_2) = $\{w_2, w_2^j, w_2^{jj}, w_2^{jjj}, \dots\}$
- Candidate (w_3) = $\{w_3, w_3^j, w_3^{jj}, w_3^{jjj}, \dots\}$

Choose the sequence W that maximizes $P(W|X)$

Noisy channel for real-world spell correction



Simplification: One error per sentence

Choose among all possible sentences with one word replaced

two of thew

- w_1, w_2^{jj}, w_3 two **off** thew
- w_1, w_2, w_3^j two of **the**
- w_1^{jjj}, w_2, w_3 **too** of thew

- So here we assume there is error in one word.
- So for assume for w_1 there are three possibilities i.e. $w_1 \rightarrow w_1' w_1''$ & w_1 (word itself)
- Similarly for $w_2 \rightarrow w_2' w_2''$ & w_2 and $w_3 \rightarrow w_3' w_3''$ & w_3
- So if we assume there is error in one word only the possible combination can be

$w_1' w_2 w_3$

$w_1'' w_2 w_3$

$w_1 w_2' w_3$

$w_1 w_2'' w_3$

$w_1 w_2 w_3'$

$w_1 w_2 w_3''$

$w_1 w_2 w_3$

- So total $6+1 = 7$ possibilities.
- Otherwise three words and three combination for each so total $3^3 = 27$.

Simplification: One error per sentence

Choose among all possible sentences with one word replaced

two of thew

- w_1, w_2, w_3 two **off** thew $w_1,$
- w_2, w_3 two of **the** w_1, w_2, w_3
- **too** of thew

Choose the sequence W that maximizes $P(W|X)$

Getting the probability values

Noisy Channel

$$\hat{W} = \operatorname{argmax}_{W \in S} P(W | X)$$

where X is the observed sentence and S is the set of all the possible sequences from the candidate set

Getting the probability values

Noisy Channel

$$\hat{W} = \operatorname{argmax}_{W \in S} P(W | X)$$

where X is the observed sentence and S is the set of all the possible sequences from the candidate set

$$= \operatorname{argmax}_{W \in S} P(X | W) P(W)$$

Getting the probability values

Noisy Channel

$$\hat{W} = \operatorname{argmax}_{W \in S} P(W | X)$$

where X is the observed sentence and S is the set of all the possible sequences from the candidate set

$$= \operatorname{argmax}_{W \in S} P(X | W) P(W)$$

$P(X|W)$

- Same as for non-word spelling correction
- Also require probability for no error $P(w|w)$

Probability of no error

What is the probability for a correctly typed word? $P(\text{"the"}|\text{"the"})$

It may depend on the source text under consideration

- 1 error in 10 words $\rightarrow 0.9$
- 1 error in 100 words $\rightarrow 0.99$

Computing $P(W)$

Use Language Model

- Unigram
- Bigram
- ...

Spell corrector

Peter Norvig

<http://www.norvig.com/spell-correct.html>

1. Spell correction algorithm based on edit distance:
<http://norvig.com/spell-correct.html>
2. Python spell correction libraries PyEnchant based on the enchant library (<http://pythonhosted.org/pyenchant/>), autocorrect, which is available at <https://github.com/phatpiglet/autocorrect/>
3. DeepSpell: Based on RNN, LSTM and word embedding
<https://github.com/MajorTal/DeepSpell>