

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](#)



Solving real-world 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
 - Task-specific classifier

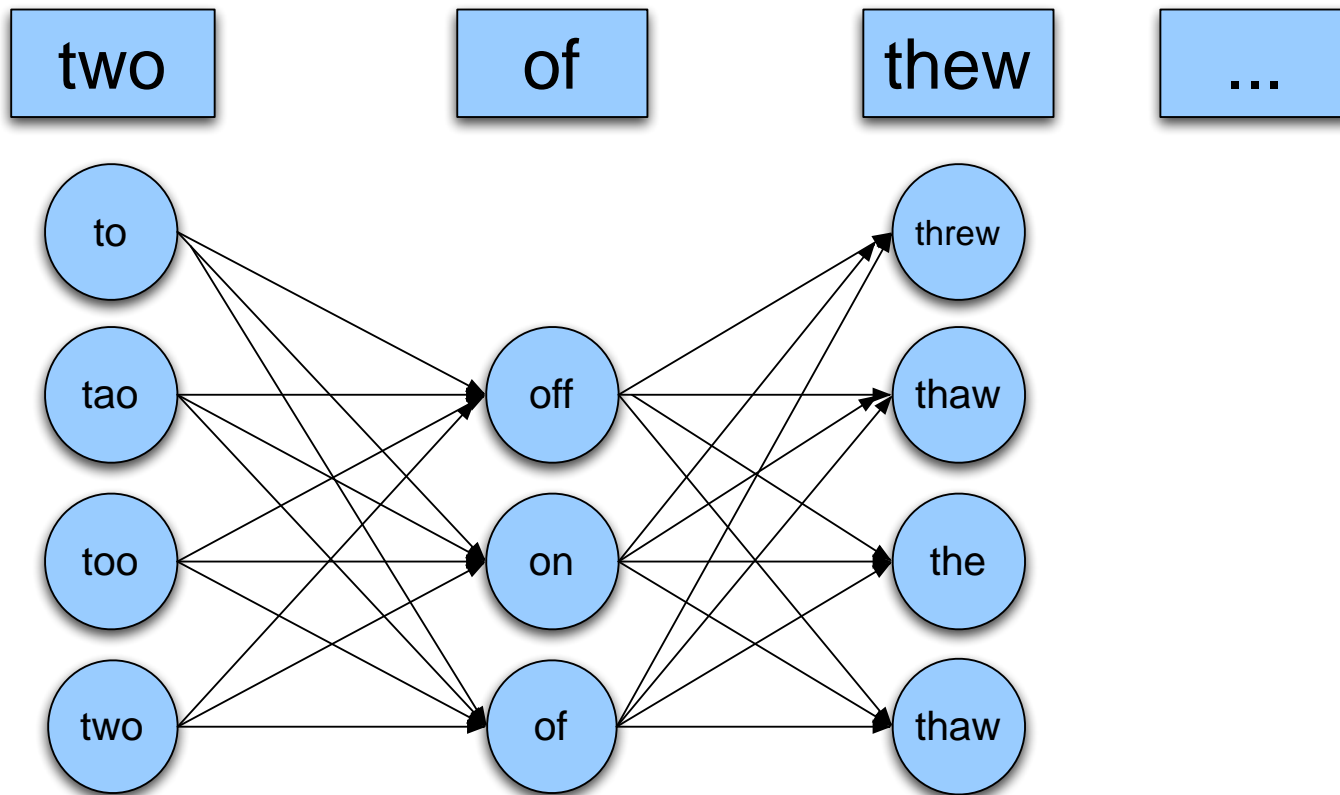


Noisy channel for real-word spell correction

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

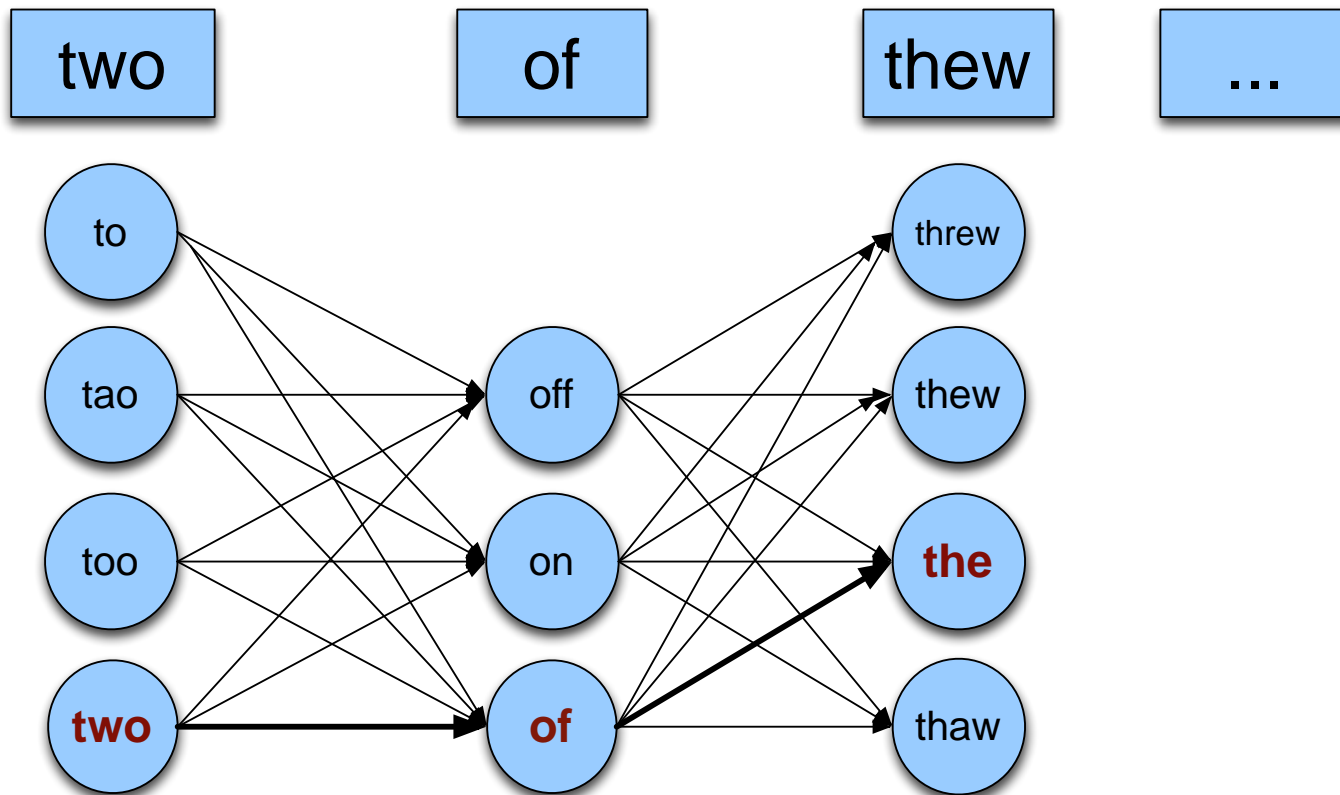


Noisy channel for real-word spell correction





Noisy channel for real-word spell correction





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 $P(W)$



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, $P(w|w)$



Probability of no error

- What is the channel probability for a correctly typed word?
- $P(\text{"the"} | \text{"the"})$
- Obviously this depends on the application
 - .90 (1 error in 10 words)
 - .95 (1 error in 20 words)
 - .99 (1 error in 100 words)
 - .995 (1 error in 200 words)



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

Spelling Correction and the Noisy Channel

Real-Word Spelling
Correction

