

←

1/28

→

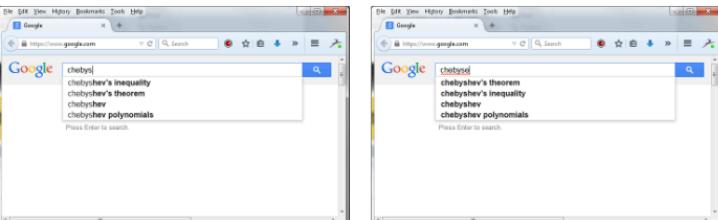
3:13:22

Auto correction/completion

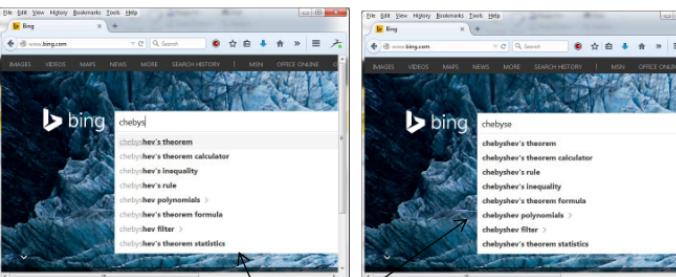
••

University of Southern California  USC

Some Google/Bing Examples

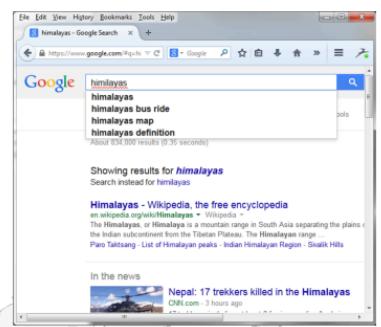


Russian mathematician, notice red underline appears as soon as the first incorrect character is typed



Bing also combines autocomplete with spelling correction but there is no red underline

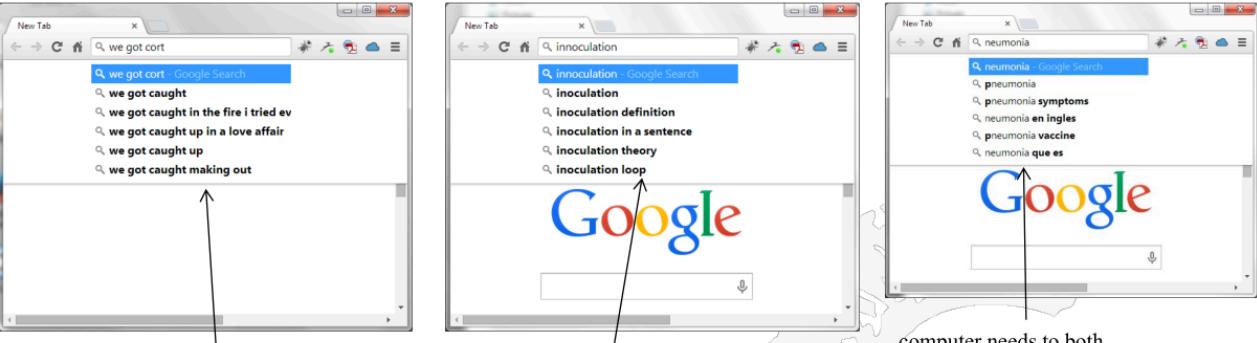
Copyright Ellis Horowitz 2011-2022



Himilayas

University of Southern California  USC

More Examples



easy for people, but harder for computers to correct, **likely use of n-grams**

easy for a computer to correct **likely use of a database** of words

computer needs to both identify the error and correct the misspelling

Google combines spelling correction with the **most likely terms** as it comes up with “cars” in autocomplete for the query “jagwa” (misspelled) but leaves the user’s misspelling for a while

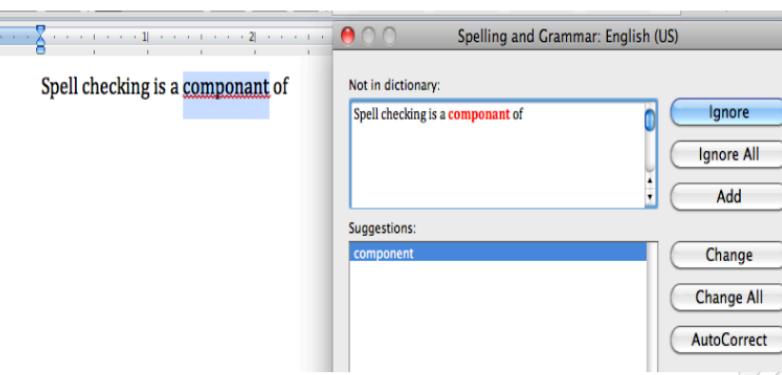
Copyright Ellis Horowitz 2011-2022

• •

University of Southern California  USC

Spelling Correction is Done in Many Places

1. Word processing



Spell checking is a component of

Not in dictionary:
Spell checking is a component of

Suggestions:
component

Ignore
Ignore All
Add
Change
Change All
AutoCorrect

2. Smartphone input



New iMessage Cancel

To: Dan Jurafsky

late x

Sorry, running layr

Send

Q W E R T Y U I O P
A S D F G H J K L
Z X C V B N M x
123 space return

Word processing is the classic application for spelling correction
Word and PowerPoint have mode to auto-correct

- set as the default
- the spell dictionary can be modified

Typing on a virtual keyboard can be doubly difficult (for seniors)

Copyright Ellis Horowitz, 2011-2022

4

..

Rates of Spelling Errors

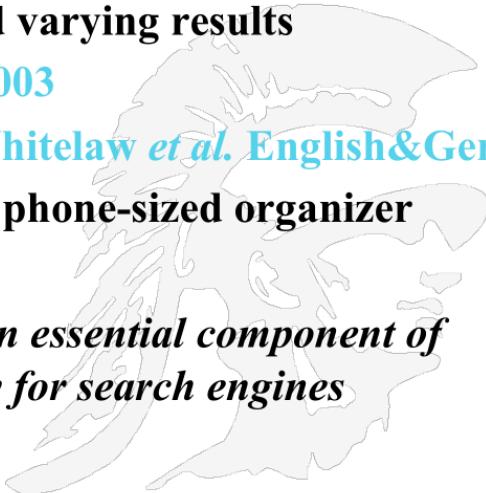
Error rates vary depending upon the application

- **Typing is very error prone, and especially difficult on smartphones**
- Different studies have produced varying results

26%: Web queries [Wang et al. 2003](#)

13%: Retyping, no backspace: [Whitelaw et al. English&German](#)

7%: Words corrected retyping on phone-sized organizer



So seamless spelling correction is an essential component of information retrieval and especially for search engines

..



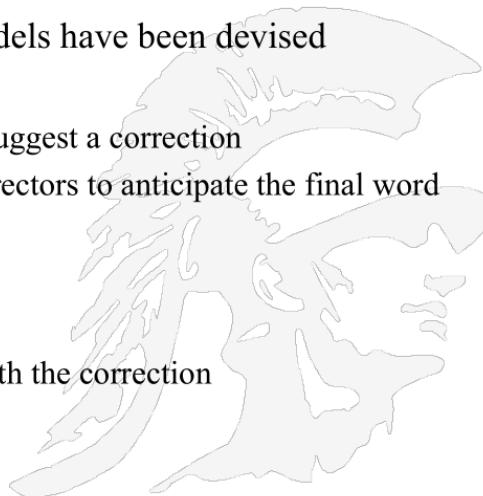
The Two Main Spelling Tasks

1. Spelling Error *Detection*

- Obviously we need a big dictionary and the ability to search it quickly
- Using context may be necessary
 - To do this spelling error models have been devised

2. Spelling Error *Correction*

- Web search engines **always** try to suggest a correction
- Autocomplete requires spelling correctors to anticipate the final word
 - Fast response time is required
- The two major techniques are
 1. edit distance algorithms or
 2. n-gram matching to come up with the correction



..

Three Types of Spelling Errors

1. Non-word errors

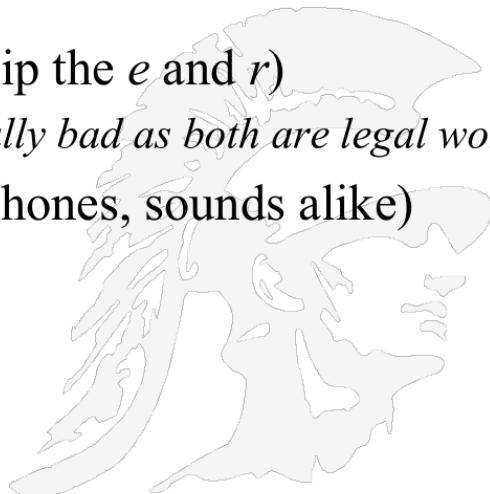
- *graffe* → *giraffe*

2. Typographical errors (flip the *e* and *r*)

- *three* → *there* (*especially bad as both are legal words*)

3. Cognitive errors (homophones, sounds alike)

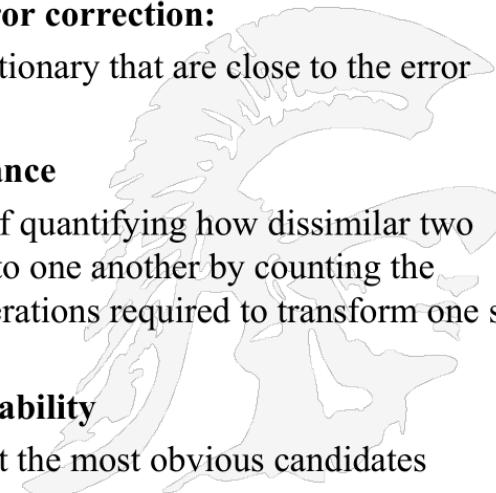
- *piece* → *peace*,
- *too* → *two*
- *your* → *you're*



..

Non-Word Spelling Errors

- **Non-word spelling error detection:**
 - Any word not in a *dictionary* is presumed to be an error
 - The larger the dictionary the better
- **Approach to non-word spelling error correction:**
 - Generate candidates from the dictionary that are close to the error
 - **How do we do this?**
 - **Shortest weighted edit distance**
 - **Edit distance** is a way of quantifying how dissimilar two strings (e.g., words) are to one another by counting the minimum number of operations required to transform one string into the other
 - **Highest noisy channel probability**
 - use probabilities to select the most obvious candidates



..



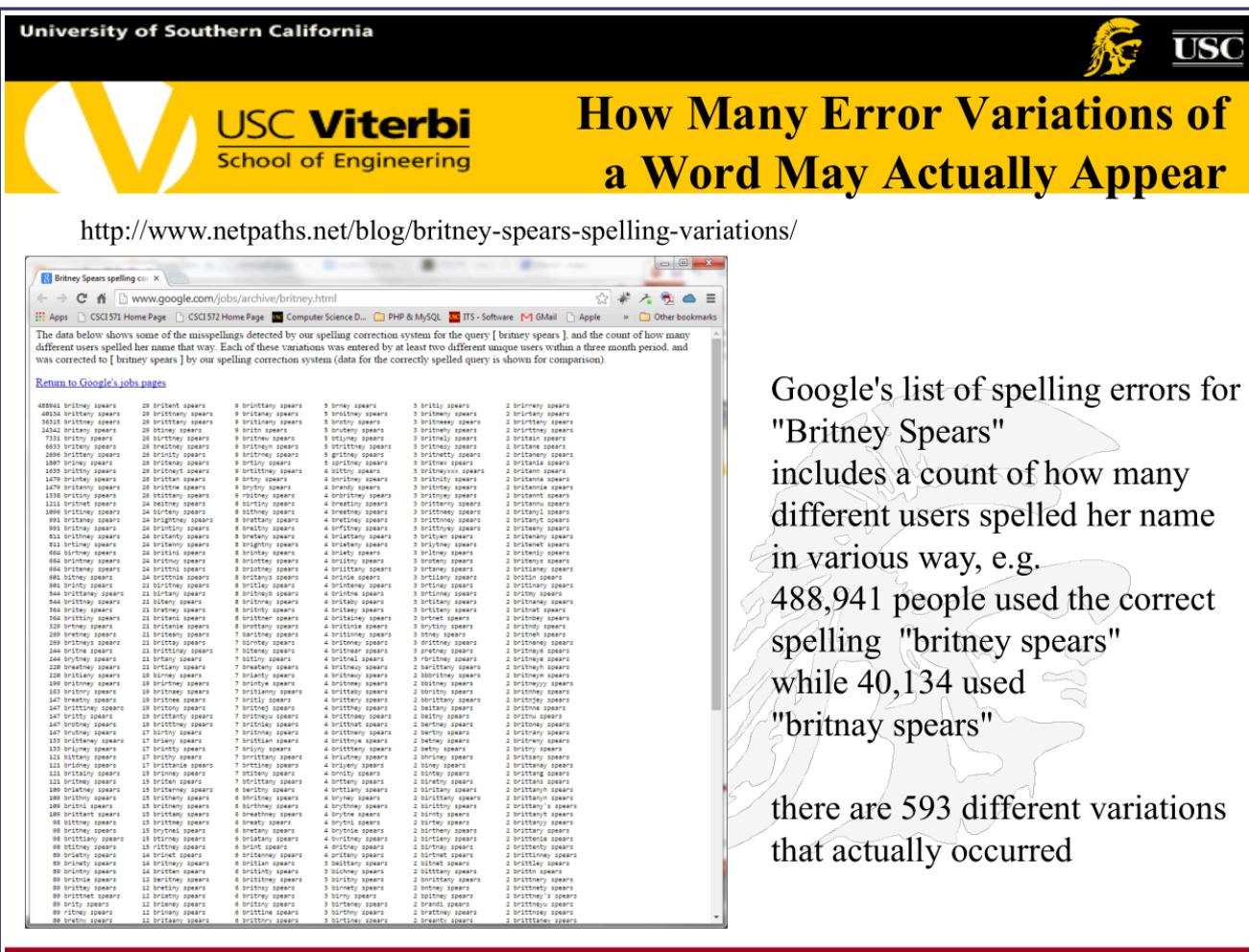
Causes of Misspellings

Cause	Misspelling	Correction
typing quickly	exxit mispell	exit misspell
keyboard adjacency	importamt	important
inconsistent rules	conceive concierge	conceive concierge
ambiguous word breaking	silver light	silverlight
new words	kinnect	kinect

According to Cucerzan and Brill, **more than 10% of search engine queries are misspelled**

"Spelling Correction as an iterative process that exploits the collective knowledge of web users"

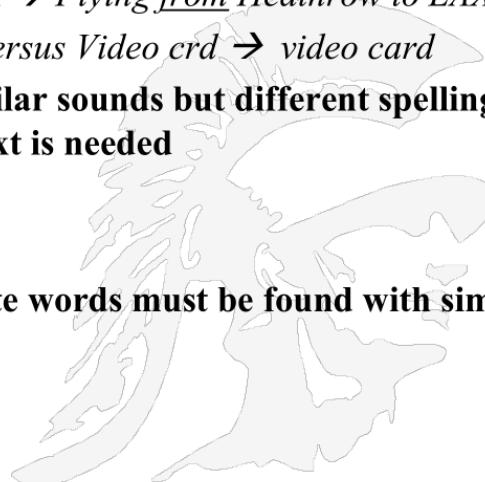
<http://csci572.com/papers/Cucerzan.pdf> (advocates using query logs to guess the correct spelling)



..

Spelling Errors Needing Context

- Some misspellings require context to disambiguate
 1. consider whether the surrounding words “make sense” for your candidate set, e.g.
 - *Flying form Heathrow to LAX → Flying from Heathrow to LAX*
 - *Power crd → power cord versus Video crd → video card*
 2. For candidate words with similar sounds but different spellings and different meanings, context is needed
 - e.g. *there, their*
 - *N-grams are most useful here*
 3. To resolve the above, candidate words must be found with similar pronunciations
 - *use the Soundex algorithm*



..

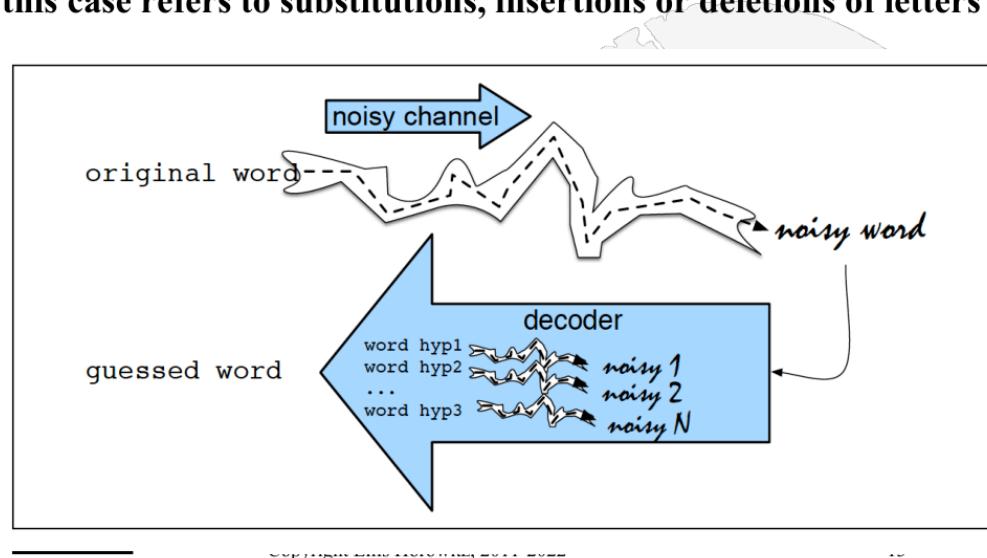
More Challenges for Identifying Spelling Errors

- Some additional challenges
 - 4. Allow for insertion of a space or hyphen
 - thisidea → this idea
 - inlaw → in-law
 - chat inspanich → chat in spanish (2 corrections)
 - 5. Allow for deletion of a space
 - power point slides → powerpoint slides
 - 6. Watch out for words NOT in the Lexicon that are valid,
e.g.
 - amd processors
 - AMD is a company that makes processors for laptops
 - Another example where context is needed

..

The Noisy Channel Model

- This model suggests treating the misspelled word as if a correctly spelled word has been distorted by being passed through a noisy communication channel
- Noise in this case refers to substitutions, insertions or deletions of letters



..

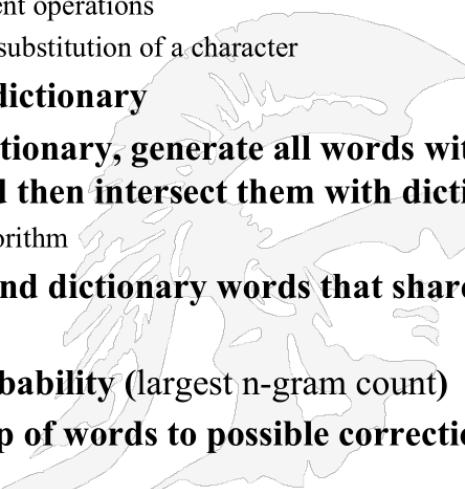
The Basic Spelling Correction Algorithm

1. **Initial step:** Create a dictionary and encode it for fast retrieval
2. When a query is submitted, the spell checker examines each word and for words not in the dictionary looks for possible character edits, namely
 - insertions,
 - deletions,
 - substitutions, and occasionally
 - transpositions
 - **Observation:**
 - 80% of errors are within edit distance 1
 - Almost all errors within edit distance 2
3. Take the output of step 2 and compute probabilities for the candidates using previously identified probability tables created from n-grams
4. Select the result with highest probability

..

The Basic Spelling Correction Algorithm Refined

- **Edit distance** is a way of quantifying how dissimilar two strings (e.g. words) are to one another by counting the minimum number of operations required to transform one string into the other
 - different algorithms assume slightly different operations
 - e.g., Levenshtein uses: removal, insertion, substitution of a character
1. **Check each query term against the dictionary**
 2. **For each term NOT found in the dictionary, generate all words within edit distance $\leq k$ (e.g., $k = 1$ or 2) and then intersect them with dictionary**
 - Compute them fast with a Levenshtein algorithm
 3. **Use a character n -gram index and find dictionary words that share “most” n -grams with word**
 4. **Select the word with the highest probability (largest n-gram count)**
 5. **For speed, have a pre-computed map of words to possible corrections**



..

Use Edit Distance To Produce Candidate Corrections

Input	Candidate Correction	Correct Letter	Error Letter	correction Type	
acress	actress	t	-	insertion	Six words within 1 of acress
acress	cress	-	a	deletion	Context check is necessary to choose the appropriate word (try this yourself in Google/Bing)
acress	caress	ca	ac	transposition	
acress	access	c	r	substitution	
acress	across	o	e	substitution	
acress	acres	-	s	deletion	

For the word "acress" there are six dictionary words all within edit distance 1

..

Now Apply Probabilities

- We now need to compute the prior probability of each occurrence
- We can do this using unigrams, bigrams, trigrams, etc
- Using the Corpus of Contemporary English, 404,253,213 words we get the following
- *Across* is the most likely choice, followed by

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

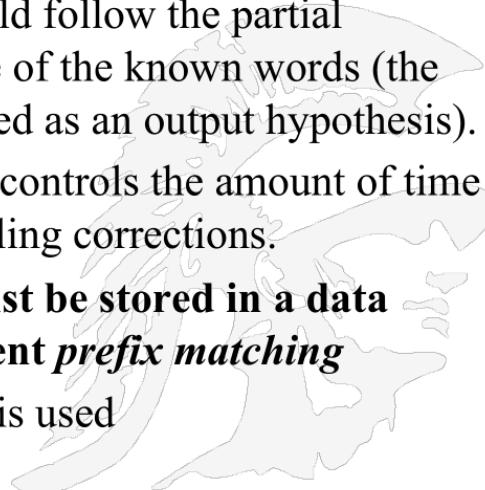
For fun try:
I need acress to ...
I love the acress ...
a kiss and a acress ...

"across" is the
most likely
correction

..

The Spelling Correction Dictionary and Autocomplete

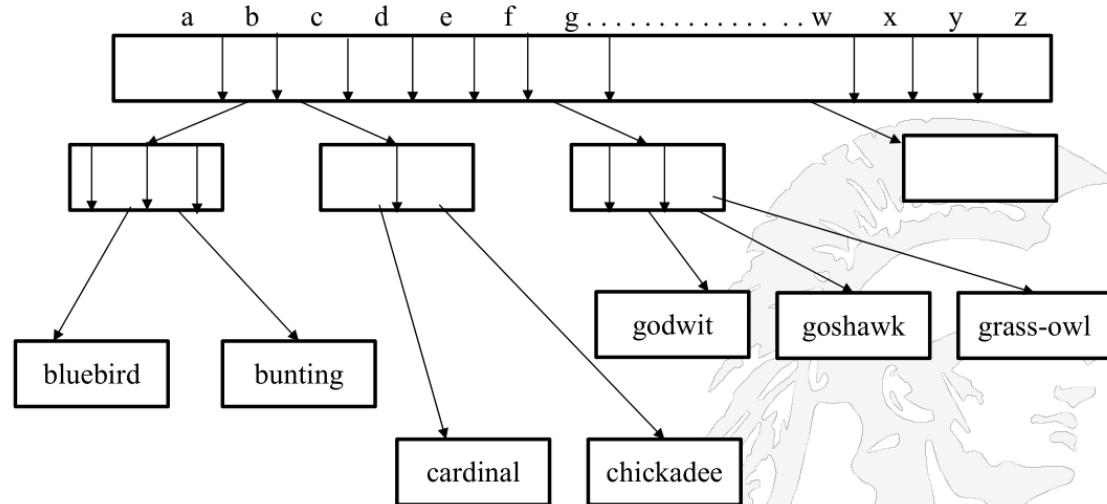
- **The search for corrections is carried out from left-to-right**
 - At each point, a partial hypothesis is expanded with every character which could follow the partial hypothesis and lead to one of the known words (the user input is always allowed as an output hypothesis).
 - Thus the branching factor controls the amount of time required to search for spelling corrections.
- **The terms of the lexicon must be stored in a data structure that affords efficient *prefix matching***
 - Often a trie data structure is used



..

The Spelling Correction Dictionary Example of a Trie

- a prefix tree (sometimes called a trie from the word retrieval) is a tree of degree ≥ 2 in which the branching at any level is determined by a portion of the key

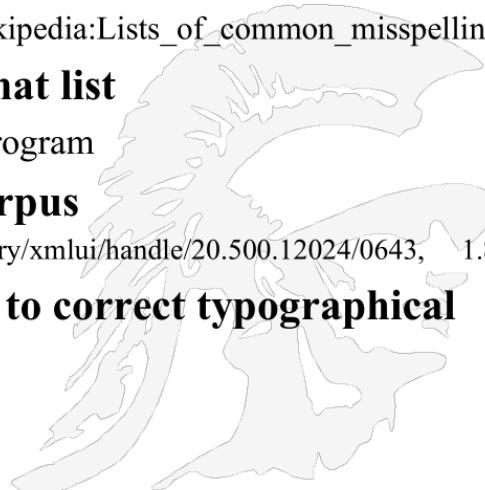


branch nodes take you down the tree to element nodes; At any stage one is pointing at all keyword matches that contain the same prefix; Computing time for retrieval is $O(m)$ where m is the length of the string, at the expense of increased storage

..

The Spelling Correction Dictionary Error Test Sets

- To enhance a lexicon one can include a table of common misspellings
- there are many possible spelling error test sets, e.g.
 - Wikipedia's list of common English misspelling
 - https://en.wikipedia.org/wiki/Wikipedia:Lists_of_common_misspellings
 - Aspell filtered version of that list
 - <http://aspell.net/> is the spell program
 - Birkbeck spelling error corpus
 - <https://ota.bodleian.ox.ac.uk/repository/xmlui/handle/20.500.12024/0643>, 1.8MBs
- These sets are primarily used to correct typographical errors



..

Using N-Grams For Spelling Correction

- An ***n*-gram model** is a type of probabilistic language model for predicting the next item in a sequence
- Two benefits of *n*-gram models (and algorithms that use them) are simplicity and scalability – with larger *n*, a model can store more context with a well-understood space–time tradeoff, enabling small experiments to scale up efficiently

- Sample of 3-gram sequences
- ceramics collectables fine (130)
- ceramics collected by (52)
- ceramics collectible pottery (50)
- ceramics collectibles cooking (45)

- Sample of 4-gram sequences and # of times appeared
- serve as the incoming (92)
 - serve as the incubator (99)
 - serve as the independent (794) ←
 - serve as the index (223)
 - serve as the indication (72)
 - serve as the indicator (120)

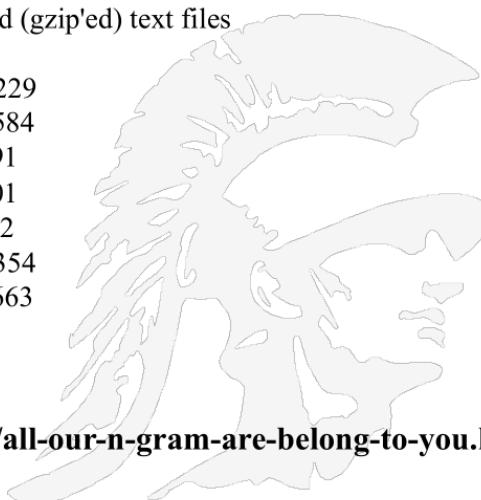
a query such as "serve as the indapendant" would match the above

Google's N-Gram Data

- Google has collected and uses a great deal of N-gram data
- Google is using the Linguistics Data Consortium to distribute more than one trillion words they have extracted from public web pages
- Below is a statistical summary of the data they are distributing

File sizes: approx. 24 GB compressed (gzip'ed) text files

Number of tokens: 1,024,908,267,229
Number of sentences: 95,119,665,584
Number of unigrams: 13,588,391
Number of bigrams: 314,843,401
Number of trigrams: 977,069,902
Number of fourgrams: 1,313,818,354
Number of fivegrams: 1,176,470,663



<http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>

Some References

- *How Difficult is it to Develop a Perfect Spell-checker? A Cross-linguistic Analysis through Complex Network Approach*, Monojit Choudhury¹, Markose Thomas², Animesh Mukherjee¹, Anupam Basu¹, and Niloy Ganguly¹

<http://citeseerx.ist.psu.edu/viewdoc/download;jsessionid=52A3B869596656C9DA285DCE83A0339F?doi=10.1.1.146.4390&rep=rep1&type=pdf>

- *Using the web for language independent spellchecking and autocorrection* by C. Whitelaw et al Proc. 2009 Conf. on Empirical Methods in Natural Language Processing, pp890-899

http://static.googleusercontent.com/external_content/untrusted_dlcp/research.google.com/en/us/pubs/archive/36180.pdf

Spell Checking by Computer, by Roger Mitton,
<http://www.dcs.bbk.ac.uk/~roger/spellchecking.html>

..

University of Southern California



Edit Distance & Levenshtein Algorithm

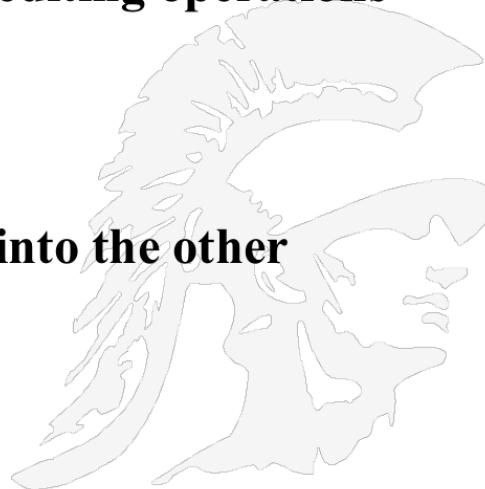


Copyright Ellis Horowitz 2011-2018

..

Edit Distance

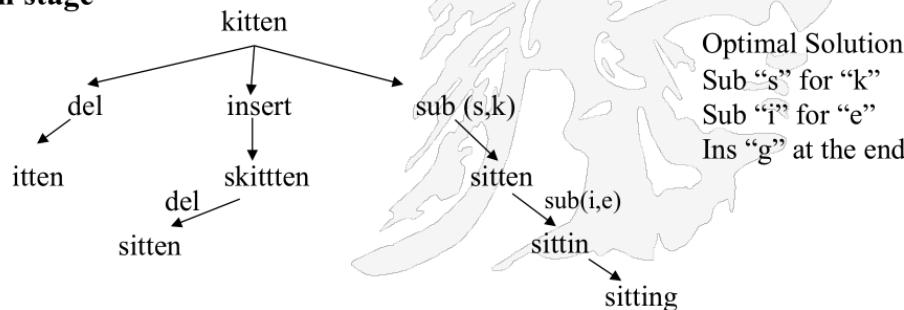
- the minimum edit distance between two strings is the minimum number of editing operations
 - insertion
 - deletion
 - substitution
- needed to transform one into the other



..

How to Find the Minimum Edit Distance

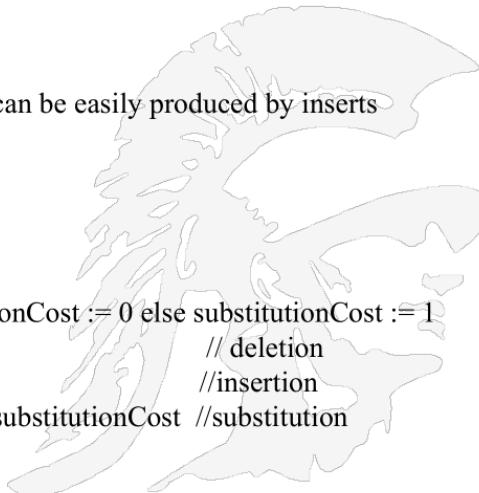
- Searching for a path (sequence of edits) from the start string to the final string
 - initial state: word we're transforming (e.g. kitten)
 - operators: insert, delete, substitute
 - goal state: the word we're trying to get to (e.g. sitting)
 - path cost: what we want to minimize, the number of edits
- If we blindly generate all possible paths in an effort to produce the goal state, our algorithm will take exponentially long
- But we realize that we needn't do that, we may only follow the path that is optimal at each stage



..

Pseudocode Implementation of Levenshtein Distance

```
function LevenshteinDistance(char s[1..m], char t[1..n]):  
    //for all i and j, d[i,j] will hold the Levenshtein distance between  
    //the first i characters of s and the first j characters of t  
    declare int d[0..m, 0..n]  
    //Set each element in d to zero  
    for i from 1 to m, j from 1 to n: d[i,j] := 0  
    Starting with empty character source and target can be easily produced by inserts  
    for i from 1 to m: d[i,0] := i  
    for j from 1 to n: d[0,j] := j  
    //main loop  
    for j from 1 to n:  
        for i from 1 to m:  
            if s[i] = t[j] then substitutionCost := 0 else substitutionCost := 1  
            d[i,j] := min (d[i-1,j] + 1, // deletion  
                           d[i,j-1] + 1, //insertion  
                           d[i-1, j-1] + substitutionCost //substitution  
            )  
    return d[m,n]
```



Copyright Ellis Horowitz 2011-2022

You can learn more about the Levenshtein algorithm [here](#).

And, you can run the algorithm [here](#), to understand how the edit distance is incrementally computed; this is a minimal version.

This is an example of how the Levenshtein distance is used in practice.



Weighted Edit Distance

- why would we add weights to the computation?
 - spell correction: some letters are more likely to be mistyped than others
- a **confusion matrix** is a specific table layout that allows visualization of the performance of an algorithm; each column of the matrix represents the instances in a predicted class while each row represents the instances in an actual class (or vice-versa)

☰ 10/102 3 - 4 - Weighted Minimum Edit Distance - Stanford NLP - Professor Dan Jurafsky & Chris Manning

Dan Jurafsky



Confusion matrix for spelling errors

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

Y (correct)

X	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
a	0	7	1	342	0	0	2	118	0	1	0	0	0	0	0	0	0	1	35	9	9	0	1	0	5	0
b	0	0	9	9	2	2	3	1	0	0	0	5	11	5	0	10	0	0	2	1	0	8	0	0	0	
c	6	5	0	16	0	9	5	0	0	0	1	0	7	9	1	10	2	5	39	40	1	3	7	1	1	0
d	1	10	13	0	12	0	5	5	0	0	2	3	7	3	0	1	0	43	30	22	0	0	4	0	2	0
e	358	0	3	11	0	2	2	0	89	0	0	3	0	5	93	0	0	14	12	6	15	0	1	0	18	0
f	0	15	0	3	1	0	2	0	0	0	0	3	0	0	0	0	0	12	0	0	2	0	0	0	0	
g	1	0	11	11	9	2	0	0	0	0	1	0	0	0	0	0	0	0	13	21	0	0	0	0	0	
h	1	8	0	3	0	0	0	0	0	0	0	2	0	12	14	2	3	0	3	11	0	0	2	0	0	
i	103	0	0	0	14	0	1	0	0	0	0	0	6	0	0	0	0	0	2	1	47	0	2	1	15	0
j	0	1	1	9	0	0	1	0	0	0	0	0	2	1	0	0	0	0	0	5	0	0	0	0	0	
k	1	2	8	4	1	1	2	5	0	0	0	0	5	0	2	0	0	0	0	6	0	0	0	4	0	0
l	1	2	10	3	4	0	4	6	13	0	0	0	0	14	2	5	0	0	10	2	0	0	0	0	0	
m	1	7	0	0	2	0	0	0	4	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0		
n	2	7	6	5	3	0	1	19	1	0	4	35	78	0	0	7	0	28	5	7	0	0	1	2	0	
o	91	1	1	3	115	0	0	0	25	0	2	0	0	0	0	14	0	2	4	14	39	0	0	18	0	
p	0	11	1	2	0	6	5	0	2	9	0	2	7	6	15	0	0	1	3	6	0	4	1	0	0	
q	0	0	1	0	0	0	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
r	0	14	0	20	12	2	0	8	4	0	0	2	4	20	1	14	0	0	12	22	0	0	1	0		
s	1	11	27	35	35	2	0	2	0	27	0	0	14	14	0	0	0	15	0	0	5	3	20	1		
t	3	4	9	42	7	5	19	5	0	1	0	14	9	5	5	6	0	11	37	0	0	2	19	0	7	
u	20	0	0	0	44	0	0	0	64	0	0	0	2	43	0	0	4	0	0	0	0	2	0	8		
v	0	0	7	0	0	3	0	0	0	0	0	1	0	0	0	0	0	0	8	3	0	0	0	0		
w	2	2	1	0	1	0	0	2	0	0	0	1	0	0	0	0	7	0	6	3	3	1	0	0		
x	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
y	0	0	2	0	15	0	1	7	15	0	0	0	2	0	6	1	0	7	36	8	5	0	0	1		
z	0	0	0	7	0	0	0	0	0	7	5	0	0	0	0	0	2	21	3	0	0	0	3			

05:33 / 2:47

the confusion matrix for spelling errors shows us, e.g. that "e" is most often confused with "a", and that "i" is often confused with both "e" and "a"