

School of Computing and Information Systems
The University of Melbourne
COMP90049 Knowledge Technologies (Semester 2, 2017)
Workshop exercises: Week 4

Suppose that we have observed the token **lended**, and we have a dictionary as follows:

addendum
blenders
commodity
deaden
end
leader
leant
lent
lemonade
pleading

1. Which, if any, of the above dictionary entries would be returned using a Neighbourhood Search with a neighbourhood of 1? 2? 3?
 - There aren't any items in the dictionary requiring only a single change from **lended**.
 - With a neighbourhood size of 2, there is a dictionary entry:
 - **leader**, by Replacing the **n** with **a**, and the second **d** with **r**
 - Along with the above, the following are also within a neighbourhood of 3:
 - **blenders**, by Inserting the **b**, Replacing the second **d** with **r**, and Inserting the **s**
 - **deaden** (three Replaces)
 - **end** (three Deletions)
 - **lent** (one Replace and two Deletions)
2. With respect to the input string **lended** and the dictionary entry **deaden**, calculate the following:
 - (a) the Global Edit Distance, using the parameter $[m, i, d, r] = [+1, -1, -1, -1]$

(a)	ε	l		e		n		d		e		d	
ε	0	\leftarrow	-1	\leftarrow	-2	\leftarrow	-3	\leftarrow	-4	\leftarrow	-5	\leftarrow	-6
	\uparrow	\nwarrow		\nwarrow		\nwarrow		\nwarrow		\nwarrow		\nwarrow	
d	-1		-1	\leftarrow	-2	\leftarrow	-3		-2	\leftarrow	-3	\leftarrow	-4
	\uparrow	\nwarrow	\uparrow	\nwarrow					\nwarrow				
e	-2		-2		0	\leftarrow	-1	\leftarrow	-2		-1	\leftarrow	-2
	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow		\nwarrow		\uparrow	\nwarrow		
a	-3		-3		-1		-1	\leftarrow	-2		-2		-2
	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow				\nwarrow	
d	-4		-4		-2		-2		0	\leftarrow	-1		-1
	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow			
e	-5		-5		-3		-3		-1		1	\leftarrow	0
	\uparrow	\nwarrow	\uparrow	\nwarrow	\uparrow	\nwarrow		\nwarrow	\uparrow		\uparrow	\nwarrow	
n	-6		-6		-4		-2		-2		0		0

- From the above table, we can observe that the Global Edit Distance is 0, corresponding to the following sequence of operations: Replace, Match, Replace, Match, Match, Replace, which I will abbreviate as **rmrmmr**. (You can follow along with the highlighted back-pointers.)

(b)	ε	l	e	n	d	e	d
ε	0	0	0	0	0	0	0
d	0	0	0	0	1	0	1
e	0	0	1	0	0	2	1
a	0	0	0	0	0	1	1
d	0	0	0	0	1	0	2
e	0	0	1	0	0	2	1
n	0	0	0	2	1	1	1

(b) the Local Edit Distance, using the parameter $[m, i, d, r] = [+1, -1, -1, -1]$

- From the above table, we can observe that the Local Edit Distance is 2 (highlighted); there are five equivalent-scoring substring matches that it corresponds to:
 - Align **-de-** in **lended** with the first **de-** in **deaden**: mm
 - Align **-ded** with **dead-**: mmim
 - Align **-de-** in **lended** with the second **-de-** in **deaden**: mm
 - Align **-ende-** with **-eade-**: mrm
 - Align **-en-** with **-en**: mm

(c) the N-Gram Distance, using $n = 2$

- We begin by generating the 2-grams of the two strings; I will use the terminal marker (#) here:
 - **lended**: #l, le, en, nd, de, ed, d#
 - **deaden**: #d, de, ea, ad, de, en, n#
- Recall that the N-Gram Distance is defined as follows:

$$D(s, t) = |G_n(s)| + |G_n(t)| - 2 \times |G_n(s) \cap G_n(t)|$$

- Here we have 7 2-grams in **lended**, as well as 7 in **deaden**. Also, the two sets share 2 2-grams: **de** and **en**. (Note that we don't double-count the **des** in **deaden**, because there is only a single **de** in **lended**)
- Consequently, the 2-gram Distance is $7 + 7 - 2 \times 2 = 10$

3. Find the best approximate match (or matches, if there are ties) in the dictionary for the string **lended**, based on the following methods; consider different parameters where necessary:

(a) the Global Edit Distance

- Using the above scoring parameter, the closest matches are **blenders** (+2) and **leader** (+2)
- You might like to try some other parameter setting, to see if they give different results.

(b) the Local Edit Distance

- Using the above scoring parameter, the closest match is **blenders** (+5)
- In this case, changing the parameter is unlikely to result in a different answer. (Why?)

(c) the N-Gram Distance

- If we are using n is 2 and padding with #, the best dictionary entry is **lent**, with a 2-Gram Distance of 6.
- You might find that removing the padding characters or changing n will give different results.

(d) Soundex

- The Soundex code of **lended** is 1533.
 - None of the dictionary entries have this exact code; however, if we permit mismatches in the Soundex code, then the best matches are **commodity** (c533), **leant** (153), **lent** (153), and **lemonade** (1553)
4. Assuming that the “correct” (intended) dictionary entry was **lent**, calculate the precision of each of the above methods of finding approximate entries from the dictionary.
- For each method, we will consider how many dictionary entries it returns as a result (predicts as a good match), as well as how many it got correct — in this case, there is only a single correct answer, so the value will be 0 or 1.
 - We have quite a few methods above!
 - Neighbourhood Search, with a neighbourhood of 1: there were any results returned from the dictionary, so precision isn’t well-defined ($\frac{0}{0}$)
 - Neighbourhood Search, with a neighbourhood of 2: there was one entry returned from the dictionary (**leader**), but it wasn’t **lent**, so the precision is $\frac{0}{1} = 0$.
 - Neighbourhood Search, with a neighbourhood of 3: there were five entries returned from the dictionary, and **lent** was one of them. The precision of this system is the number of correct responses (1) out of the total number of attempted responses (5), $\frac{1}{5} = 20\%$
 - Global Edit Distance: there were two results from the dictionary (**blenders** and **leader**), but no **lent**, so the precision is $\frac{0}{2} = 0$
 - Local Edit Distance: there was just a single result (**blenders**) which wasn’t **lent**, so the precision is 0
 - N-Gram Distance: there was a single result which was **lent**, so the precision is $\frac{1}{1} = 100\%$
 - Soundex: there weren’t any exact matches with the Soundex code of **lended**, so precision isn’t well defined
 - Soundex (allowing approximate matches): allowing approximate matches of the Soundex code meant that there were four results, including **lent**, so the precision is $\frac{1}{4} = 25\%$
 - Here, the best method according to precision was the N-Gram Distance. However, if we wanted to seriously compare these methods, we would need to aggregate these results over a large number of inputs; just considering a single word is not enough information to draw solid conclusions.
 - (In fact, the fact that the N-Gram Distance uniquely found the correct result was mostly an accident of its bias toward shorter strings; you might like to think about why this happens.)